# MIRROR DESCENT-ASCENT FOR MEAN-FIELD MIN-MAX PROBLEMS

# **Anonymous authors**

Paper under double-blind review

# **ABSTRACT**

We study two variants of the mirror descent-ascent algorithm for solving minmax problems on the space of measures: simultaneous and sequential. We work under assumptions of convexity-concavity and relative smoothness of the payoff function with respect to a suitable Bregman divergence, defined on the space of measures via flat derivatives. We show that the convergence rates to mixed Nash equilibria, measured in the Nikaidò-Isoda error, are of order  $\mathcal{O}\left(N^{-1/2}\right)$  and  $\mathcal{O}\left(N^{-2/3}\right)$  for the simultaneous and sequential schemes, respectively, which is in line with the state-of-the-art results for related finite-dimensional algorithms.

# 1 Introduction

Numerous tasks in machine learning can be framed as optimization problems for functions defined on the space of probability measures. For instance, in supervised learning, pioneering works (Chizat & Bach, 2018; Mei et al., 2018; Rotskoff & Vanden-Eijnden, 2018) showed that training a shallow neural network (NN) in the mean-field regime (i.e., an infinite-width one-hidden-layer NN) can be viewed as minimizing a convex function over the space of probability distributions of the parameters of the network. This key insight proved to be a fruitful approach in analyzing convergence of training algorithms for infinite-width one-hidden-layer NNs (see, e.g., (Hu et al., 2021; Chizat, 2022a; Nitanda et al., 2022; Suzuki et al., 2023)).

The paradigm of mean-field optimization has been extended to min-max settings in several works, e.g., (Hsieh et al., 2019; Domingo-Enrich et al., 2020; Wang & Chizat, 2023; Lu, 2023; Trillos & Trillos, 2023; Kim et al., 2024), which formulate the training of Generative Adversarial Networks (GANs) and adversarial robustness as a problem of finding mixed Nash equilibria (MNEs) of min-max games over the space of probability measures.

In this work, we study the convergence of an infinite-dimensional mirror descent-ascent algorithm (MDA) to mixed Nash equilibria of a min-max game with a convex-concave payoff function. In games, the design of learning algorithms heavily depends on the playing conventions the players can adopt: simultaneous (players move at the same time) or sequential (each player moves upon observing the opponents' moves). To our knowledge, the works concerned with studying the convergence of discrete-time algorithms for mean-field min-max games only analyze the case of simultaneous playing (see, e.g., (Hsieh et al., 2019; Wang & Chizat, 2023)). In contrast, we make a rigorous comparison between the simultaneous and sequential algorithms, and prove that sequential playing leads to faster convergence rate. This result theoretically underpins the common practice of training GANs in an alternating fashion.

#### 1.1 NOTATION AND SETUP

For any  $\mathcal{X} \subset \mathbb{R}^d$ , let  $\mathcal{P}(\mathcal{X})$  denote the set of probability measures on  $\mathcal{X}$ . In game theory, if  $\mathcal{X}$  is the set of *(pure) strategies* available to the players, then  $\mathcal{P}(\mathcal{X})$  is known as the set of *mixed strategies*. Let  $\mathcal{C}, \mathcal{D} \subseteq \mathcal{P}(\mathcal{X})$  be convex. We consider a convex-concave (cf. Assumption 1.5) payoff function  $F: \mathcal{C} \times \mathcal{D} \to \mathbb{R}$  and the associated min-max game

$$\min_{\nu \in \mathcal{C}} \max_{\mu \in \mathcal{D}} F(\nu, \mu). \tag{1}$$

We are interested in finding *mixed Nash equilibria* (MNEs) for game (1), i.e., pairs of strategies  $(\nu^*, \mu^*) \in \mathcal{C} \times \mathcal{D}$  such that, for any  $(\nu, \mu) \in \mathcal{C} \times \mathcal{D}$ , we have

$$F(\nu^*, \mu) \le F(\nu^*, \mu^*) \le F(\nu, \mu^*).$$
 (2)

We observe that in the case in which F is bilinear, i.e.,  $F(\nu,\mu)=\int_{\mathcal{X}}\int_{\mathcal{X}}f(x,y)\nu(\mathrm{d}x)\mu(\mathrm{d}y)$ , for some  $f:\mathcal{X}\times\mathcal{X}\to\mathbb{R}$ , measures characterized by (2) are MNEs in the classical sense of two-player zero-sum games. Throughout, we assume that there exists at least one MNE for game (1).

In min-max games, the distance between a pair of strategies  $(\nu, \mu)$  and an MNE is typically measured using the Nikaidò-Isoda (NI) error (Nikaidô & Isoda, 1955), which, for all  $(\nu, \mu) \in \mathcal{C} \times \mathcal{D}$ , is defined by

$$\mathrm{NI}(\nu,\mu) \coloneqq \max_{\mu' \in \mathcal{D}} F(\nu,\mu') - \min_{\nu' \in \mathcal{C}} F(\nu',\mu).$$

Straight from the definition, we see that  $NI(\nu, \mu) \ge 0$  for all  $(\nu, \mu) \in \mathcal{C} \times \mathcal{D}$ , and from (2) it follows that  $NI(\nu, \mu) = 0$  if and only if  $(\nu, \mu)$  is an MNE.

# 1.2 MOTIVATING EXAMPLE: TRAINING OF GANS

Let  $\hat{\xi} \in \mathcal{P}(\mathcal{Y})$  be the empirical measure of the i.i.d. sampled particles  $\{x_i\}_{i=1}^M \subset \mathcal{Y}$ , and let  $\xi \in \mathcal{P}(\mathcal{Z})$  be a source measure. Consider the measurable parametrized transport map  $T_\theta : \mathcal{Z} \to \mathcal{Y}$  (which typically can be viewed as a neural network with parameters  $\theta \in \Theta \subset \mathbb{R}^d$ ). The *pushforward* of the measure  $\xi$  on  $\mathcal{Z}$  via  $T_\theta$  is the measure  $T_\theta \# \xi$  on  $\mathcal{Y}$  characterized by  $\int_{\mathcal{Y}} \varphi d(T_\theta \# \xi) = \int_{\mathcal{Z}} (\varphi \circ T_\theta) d\xi$ , for any measurable function  $\varphi : \mathcal{Y} \to \mathbb{R}$ .

The aim of a GAN is to search for the optimal set of parameters  $\theta^* \in \Theta$  that minimizes the distance between the generated measure  $T_{\theta^*} \# \xi$  and the empirical measure  $\hat{\xi}$ . In order to evaluate this distance, we define the function  $D_w: \mathcal{Y} \to \mathbb{R}$  (which can also be viewed as a neural network with parameters  $w \in \mathcal{W} \subset \mathbb{R}^d$ ), and solve the min-max problem

$$\min_{\theta \in \Theta} \max_{w \in \mathcal{W}} \left\{ \int_{\mathcal{Y}} D_w(y) \left( T_{\theta} \# \xi - \hat{\xi} \right) (\mathrm{d}y) \right\}.$$

For example, if the family of functions  $\{D_w\}_{w\in\mathcal{W}}$  is either 1-Lipschitz continuous or uniformly bounded, the resulting GAN corresponds to the Wasserstein GAN or the Total Variation GAN, respectively (Arjovsky et al., 2017). On the other hand, if the family of functions  $\{D_w\}_{w\in\mathcal{W}}$  belongs to the norm unit ball of a reproducing kernel Hilbert space (RKHS), we recover the Maximum Mean Discrepancy (MMD) GAN (Li et al., 2017).

Solving this problem on the finite-dimensional subspaces  $\theta, w \subset \mathbb{R}^d$  may pose serious challenges such as the lack of existence of pure Nash equilibria. Instead, we lift the problem to the space of probability measures and search for MNEs, i.e., optimal distributions over the set of parameters.

That is, by setting  $f(\theta, w) := \int_{\mathcal{Y}} D_w(y) \left( T_{\theta} \# \xi - \hat{\xi} \right) (\mathrm{d}y)$ , we solve the mean-field min-max game

$$\min_{\nu \in \mathcal{P}(\Theta)} \max_{\mu \in \mathcal{P}(\mathcal{W})} \left\{ \int_{\mathcal{W}} \int_{\Theta} f(\theta, w) \nu(\mathrm{d}\theta) \mu(\mathrm{d}w) \right\}. \tag{3}$$

We will demonstrate theoretically (cf. Theorem 2.4 and Theorem 3.6) that sequential updates speed up GANs training significantly. Note that the lifted problem is bilinear in  $\nu$  and  $\mu$ , so an MNE for (3) exists under mild conditions (see footnote 1).

We stress, however, that our framework applies more broadly, and, while encompassing (3) as a special case, covers also more general nonlinear convex-concave functions F. An example of an application where a nonlinear F arises naturally is discussed in Example E.

#### 1.3 RELATED WORKS

Mirror descent (MD) was originally proposed in (Nemirovski & Yudin, 1983) for solving convex optimization problems and has been extensively studied on finite-dimensional vector spaces, see e.g.

 $<sup>^1</sup>$  If F is continuous and  $\mathcal D$  is compact, then the existence of an MNE of (1) follows from Sion's minimax theorem (Sion, 1958). For the particular case when  $F(\nu,\mu)=\int_{\mathcal X}\int_{\mathcal X}f(x,y)\nu(\mathrm{d}x)\mu(\mathrm{d}y)$ , an MNE exists due to Glicksberg's minimax theorem (Glicksberg, 1952) if f is continuous and  $\mathcal C,\mathcal D$  are compact.

(Beck & Teboulle, 2003; Bubeck, 2015; Lu et al., 2018). One of its main advantages over traditional gradient descent is that, by utilizing Bregman divergence as a regularization term instead of the usual squared Euclidean norm, the MD method captures the geometry of the ambient space better than the gradient descent scheme (see (Beck & Teboulle, 2003) for a detailed discussion).

Recently, the MD algorithm has been extended to infinite-dimensional settings for studying optimization problems over spaces of measures, with applications in machine learning (e.g., Sinkhorn's and Expectation–Maximization algorithms, see (Aubin-Frankowski et al., 2022)) as well as in policy optimization for reinforcement learning (Tomar et al., 2021; Kerimkulov et al., 2023).

By leveraging results from optimization on  $\mathbb{R}^d$  (see (Bauschke et al., 2017; Lu et al., 2018)), the work of (Aubin-Frankowski et al., 2022) extends the convergence proof from (Lu et al., 2018) to the case of the infinite-dimensional MD method by showing that in order for the MD procedure to converge with rate  $\mathcal{O}\left(N^{-1}\right)$ , it suffices to require convexity and relative smoothness of F (cf. Assumptions 1.5 and 1.6, respectively).

Other works such as (Hsieh et al., 2019; Dvurechensky & Zhu, 2024) studied infinite-dimensional MDA and Mirror Prox algorithms for finding MNEs of two-player zero-sum games. The most closely related work to ours is (Hsieh et al., 2019), which focuses on min-max games for bilinear objective functions and utilizes a particular case of the MDA algorithm with relative entropy regularization.

Our paper generalizes the setting of (Hsieh et al., 2019) by considering a possibly non-linear convex-concave objective function and the MDA algorithm with a general Bregman divergence. Moreover, while (Hsieh et al., 2019) proves an explicit convergence rate  $\mathcal{O}\left(N^{-1/2}\right)$  only for the simultaneous MDA algorithm, we also prove a faster convergence rate  $\mathcal{O}\left(N^{-2/3}\right)$  for the sequential scheme. For a brief discussion on recent results on related Mirror Prox algorithms (not studied in the present paper), see Appendix J.

#### 1.4 OUR CONTRIBUTION

We provide a theoretical analysis of the proposed simultaneous and sequential MDA algorithms, establishing convergence rates under convexity–concavity and relative smoothness of the objective F with respect to a Bregman divergence. In particular, Theorem 2.4 and 3.6 show that the sequential MDA scheme achieves faster convergence than the simultaneous one. We validate our results on simple numerical experiments.

From one perspective, our work extends (Aubin-Frankowski et al., 2022) to the setting of min-max games. A key obstacle we overcome is that, unlike in single-player MD in both infinite-dimensional (cf. (Aubin-Frankowski et al., 2022)) and finite-dimensional (cf. (Lu et al., 2018)) settings, the objective function for min-max problems is not monotonically decreasing along the iterates, which forces us to work with the NI error and requires different proof techniques.

From another perspective, we generalize the results of (Hsieh et al., 2019) by considering a possibly non-linear convex—concave objective function and MDA algorithms with respect to a general Bregman divergence. Whereas (Hsieh et al., 2019) derive an explicit convergence rate only for the simultaneous MDA algorithm in the context of GAN training, we establish a faster rate for the sequential variant. Moreover, our more general framework also covers applications other than GANs, such as adversarial training of neural networks (see Example E).

At the technical level, our convergence proof for sequential MDA relies on a duality between the Bregman divergence on the space of measures and a corresponding dual Bregman divergence defined on the space of bounded measurable functions. To our knowledge, the use of this dual formulation of MDA on a function space is novel, and may be of independent interest.

# 1.5 Bregman divergence on the space of probability measures

As noted in Section 1.3, the MD algorithm relies on the Bregman divergence. We now introduce this concept rigorously for the space of probability measures using the flat derivative (Definition F.1), following (Aubin-Frankowski et al., 2022), who defined it via directional derivatives.

Set  $\mathcal{E} := \mathcal{C} \cup \mathcal{D} \subseteq \mathcal{P}(\mathcal{X})$  and let  $h : \mathcal{P}(\mathcal{X}) \to \mathbb{R}$  satisfy the following assumption.

**Assumption 1.1** (Differentiability and convexity of h). Assume that h is lower semi-continuous on  $\mathcal{E}$  and admits first-order flat derivative (cf. Definition F.1) on  $\mathcal{E}$ . Moreover, assume that  $h: \mathcal{E} \to \mathbb{R}$  is strictly convex on  $\mathcal{E}$ , i.e., for all  $\lambda \in [0,1]$  and all  $\nu', \nu \in \mathcal{E}$ , we have  $h((1-\lambda)\nu + \lambda\nu') < (1-\lambda)h(\nu) + \lambda h(\nu')$ .

If Assumption 1.1 holds, then it can be shown, via (Hu et al., 2021, Lemma 4.1), that h is strictly convex on  $\mathcal{E}$  in the sense of Assumption 1.5, i.e., for any  $\nu'$ ,  $\nu \in \mathcal{E}$ , we have

$$h(\nu') - h(\nu) > \int_{\mathcal{X}} \frac{\delta h}{\delta \nu} (\nu, x) (\nu' - \nu) (\mathrm{d}x).$$

Under Assumption 1.1, we define the h-Bregman divergence (or simply Bregman divergence) on the space of probability measures.

**Definition 1.2** (Bregman divergence). The h-Bregman divergence is the map  $D_h : \mathcal{E} \times \mathcal{E} \to [0, \infty)$  given by

$$D_h(\nu',\nu) := h(\nu') - h(\nu) - \int_{\mathcal{X}} \frac{\delta h}{\delta \nu} (\nu, x) (\nu' - \nu) (\mathrm{d}x).$$

We observe that, by Assumption 1.1,  $\int_{\mathcal{X}} \frac{\delta h}{\delta \nu}(\nu,x)(\nu'-\nu)(\mathrm{d}x)$  is well-defined, and that  $D_h(\nu',\nu) \geq 0$ , for all  $\nu',\nu\in\mathcal{E}$ , with equality if and only if  $\nu'=\nu$ .

We now give two examples of a function h and the corresponding sets  $\mathcal E$  such that Assumption 1.1 is satisfied.

**Example 1.3** (Relative entropy). Suppose that h is the relative entropy, i.e.,  $h(\nu) := \int_{\mathcal{X}} \log \frac{\nu(x)}{\pi(x)} \nu(x) dx$ , where  $\nu, \pi \in \mathcal{P}_{\lambda}(\mathcal{X})$ , i.e., they are absolutely continuous with respect to the Lebesgue measure on  $\mathcal{X}$  and  $\pi$  is a fixed reference probability measure on  $\mathcal{P}_{\lambda}(\mathcal{X})$ . Fix  $\beta > 0$ 

and define  $\mathcal{E}_{\beta} := \left\{ \nu \in \mathcal{P}_{\pi}(\mathcal{X}) : \left\| \log \frac{\nu(\cdot)}{\pi(\cdot)} \right\|_{L^{\infty}(\mathcal{X})} \le \beta \right\}$ . Note that  $\mathcal{E}_{\beta}$  is convex. From (Dupuis &

Ellis, 1997, Lemma 1.4.3), we know that the relative entropy is lower semi-continuous on  $\mathcal{P}(\mathcal{X})$ , hence on  $\mathcal{E}_{\beta}$ . Clearly, we see that h is strictly convex on  $\mathcal{E}_{\beta}$  due to the strict convexity of the map  $(0,\infty)\ni z\mapsto z\log z$ . Moreover, it is proved in (Kerimkulov et al., 2024, Proposition 2.16) that h admits the flat derivative

$$\frac{\delta h}{\delta \nu}(\nu, x) = \log \frac{\nu(x)}{\pi(x)} - h(\nu), \tag{4}$$

on  $\mathcal{E}_{\beta}$ , and for all  $\nu, \nu' \in \mathcal{E}_{\beta}$ , the Bregman divergence  $D_h(\nu', \nu)$  is in fact the Kullback-Leibler divergence (or relative entropy)  $\mathrm{KL}(\nu', \nu)$ .

**Example 1.4** ( $\chi^2$ -divergence). Suppose that h is the  $\chi^2$ -divergence, i.e.,  $h(\nu) := \frac{1}{2} \int_{\mathcal{X}} \left( \frac{\nu(x)}{\pi(x)} - 1 \right)^2 \pi(x) \mathrm{d}x$ , where  $\nu, \pi \in \mathcal{P}_{\lambda}(\mathcal{X})$ . Let  $L^2_{\pi}(\mathcal{X})$  be the set of square integrable func-

tions on 
$$\mathcal{X}$$
 with respect to  $\pi$ . Fix  $\eta > 0$  and define  $\mathcal{F}_{\eta} \coloneqq \left\{ \nu \in \mathcal{P}_{\pi}(\mathcal{X}) : \left\| \frac{\nu(\cdot)}{\pi(\cdot)} \right\|_{L^{2}_{\pi}(\mathcal{X})} \leq \eta \right\}$ . Note

that  $\mathcal{F}_{\eta}$  is convex. From (Ambrosio et al., 2000, Theorem 2.34), we know that the the  $\chi^2$ -divergence is lower semi-continuous on  $\mathcal{P}(\mathcal{X})$ , hence on  $\mathcal{F}_{\eta}$ . Clearly, we see that h is strictly convex on  $\mathcal{F}_{\eta}$  due to the strict convexity of the map  $(0,\infty)\ni z\mapsto (z-1)^2$ . Moreover, it is proved in (Kerimkulov et al., 2024, Proposition 2.18) that h admits the flat derivative  $\frac{\delta h}{\delta \nu}(\nu,x)=\frac{\nu(x)}{\pi(x)}-\int_{\mathbb{R}^d}\frac{\nu(x)}{\pi(x)}\nu(x)\mathrm{d}x$ , on  $\mathcal{F}_{\eta}$ , and for all  $\nu,\nu'\in\mathcal{F}_{\eta}$ , the Bregman divergence  $D_h(\nu',\nu)$  is in fact the  $L^2$ -distance  $\frac{1}{2}\left\|\frac{\nu'(\cdot)}{\pi(\cdot)}-\frac{\nu(\cdot)}{\pi(\cdot)}\right\|_{L^2(\mathcal{X})}^2$ .

For other examples of regularizers h that verify Assumption 1.1 and frequently appear in machine learning applications, see (Kerimkulov et al., 2024, Proposition 2.20).

#### 1.6 SIMULTANEOUS AND SEQUENTIAL MDA

In what follows, we state our standing assumptions and the necessary definitions for introducing the simultaneous and sequential MDA schemes. Let  $F: \mathcal{C} \times \mathcal{D} \to \mathbb{R}$  be such that  $\nu \mapsto F(\nu, \mu)$  and  $\mu \mapsto F(\nu, \mu)$  admit first-order flat derivatives (cf. Definition F.1) on  $\mathcal{C}$  and  $\mathcal{D}$ , respectively.

**Assumption 1.5** (Convexity-concavity of F). Assume that F is convex in  $\nu$  and concave in  $\mu$ , i.e., for any  $\nu, \nu' \in \mathcal{C}$  and any  $\mu, \mu' \in \mathcal{D}$ , we have  $D_{F(\cdot, \mu)}(\nu', \nu) \geq 0$  and  $D_{F(\nu, \cdot)}(\mu', \mu) \leq 0$ .

**Assumption 1.6** (Relative smoothness of F). Assume that, given  $L_{\nu}, L_{\mu} > 0$ , the function F is  $L_{\nu}$ -smooth in  $\nu$  and  $L_{\mu}$ -smooth in  $\mu$  relative to h, i.e., for any  $\nu, \nu' \in \mathcal{C}$  and any  $\mu, \mu' \in \mathcal{D}$ , we have  $D_{F(\cdot,\mu)}(\nu',\nu) \leq L_{\nu}D_h(\nu',\nu)$  and  $D_{F(\nu,\cdot)}(\mu',\mu) \geq -L_{\mu}D_h(\mu',\mu)$ .

In Proposition D.1 and E.1, we verify that Assumption 1.5 and 1.6 are satisfied by Example 1.2 and E. In Lemma C.3, we show that Assumption 1.5 and 1.6 correspond to the intuition we have from optimization on  $\mathbb{R}^d$ , where convexity and relative smoothness are equivalent, respectively, to the Hessian of F being non-negative, and upper bounded by the Hessian of h weighted by the smoothness constant.

For a given stepsize  $\tau > 0$ , and fixed initial pair of strategies  $(\nu_0, \mu_0) \in \mathcal{C} \times \mathcal{D}$ , for  $n \geq 0$ , the *simultaneous* and *sequential* MDA algorithms are respectively defined by

# Algorithm 1: SIMULTANEOUS MDA

#### **Algorithm 2: SEQUENTIAL MDA**

```
Input: Objective function F, initial measures (\nu^0, \mu^0), stepsize \tau > 0 for n = 0, 1, \dots, N-1 do  \begin{bmatrix} \nu^{n+1} = \arg\min_{\nu \in \mathcal{C}} \{\int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^n, \mu^n, x) (\nu - \nu^n) (\mathrm{d}x) + \frac{1}{\tau} D_h(\nu, \nu^n) \}, \\ \mu^{n+1} = \arg\max_{\mu \in \mathcal{D}} \{\int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n+1}, \mu^n, y) (\mu - \mu^n) (\mathrm{d}y) - \frac{1}{\tau} D_h(\mu, \mu^n) \} \end{bmatrix} Output: \left(\frac{1}{N} \sum_{n=0}^{N-1} \nu^{n+1}, \frac{1}{N} \sum_{n=0}^{N-1} \mu^n \right)
```

Although we abuse the notation by denoting both (1) and (2) by  $(\nu^n, \mu^n)_{n\geq 0}$ , we will make it clear from the context which algorithm we consider.

Algorithm (1) is referred to as *simultaneous* because both players update their strategy from step n to n+1 at the same time, whereas Algorithm (2) is called *sequential* because the minimizing player is first updating their move from step n to n+1, and then the maximizing player is acting upon observing the minimizing player's (n+1)-th action. Note that due to the symmetry of the players, the analysis of scheme (2) also covers the case when the maximizing player moves first followed by the minimizing player.

The motivation behind the use of the terms involving  $\frac{\delta F}{\delta \nu}$  and  $\frac{\delta F}{\delta \mu}$  in algorithms (1) and (2), is that instead of minimizing and maximizing directly on F (which could be a potentially intractable problem), we minimize and maximize over  $\nu$  and  $\mu$  in the first-order linear approximations  $F(\nu^n,\mu^n)+\int_{\mathcal{X}}\frac{\delta F}{\delta \nu}(\nu^n,\mu^n,x)(\nu-\nu^n)(\mathrm{d}x)$  and  $F(\nu^n,\mu^n)+\int_{\mathcal{X}}\frac{\delta F}{\delta \mu}(\nu^n,\mu^n,y)(\mu-\mu^n)(\mathrm{d}y)$ . In order to make sure that these approximations around  $(\nu^n,\mu^n)$  are precise enough, we penalize the distance between  $(\nu^{n+1},\mu^{n+1})$  and  $(\nu^n,\mu^n)$  by introducing the Bregman regularization terms  $\frac{1}{\tau}D_h(\nu,\nu^n)$  and  $\frac{1}{\tau}D_h(\mu,\mu^n)$ .

We observe that by varying the choices of h in Definition 1.2 we obtain a collection of different update rules in the MDA algorithms (1) and (2). When h is the relative entropy, we can view (1) and (2) as Euler discretizations of a Fisher-Rao gradient flow, whose continuous-time convergence with

explicit rates for mean-field min-max games was proved in (Lascu et al., 2024) (cf. also (Liu et al., 2023) for single-player convex optimization).

# 2 CONVERGENCE OF THE SIMULTANEOUS MDA ALGORITHM (1)

In this section, we state the main result on the convergence of the simultaneous MDA algorithm. Proving that the simultaneous MDA algorithm (1) converges relies on the following key assumption.

**Assumption 2.1.** Suppose that F is Lipschitz relative to h, i.e., there exists  $L_F > 0$  such that, for any  $\nu, \nu' \in C$  and any  $\mu, \mu' \in D$ ,

$$|F(\nu', \mu') - F(\nu, \mu)|^2 \le L_F (D_h(\nu', \nu) + D_h(\mu', \mu)).$$

**Remark 2.2.** We show in Lemma C.2 that Assumption 2.1 is satisfied (via Pinsker's inequality, that is,  $TV^2(\nu',\nu) \leq \frac{1}{2} KL(\nu',\nu)$ ) when F has bounded first-order flat derivatives in  $\nu$  and  $\mu$ , and h is the relative entropy, i.e.,  $h(\nu) \coloneqq \int_{\mathcal{X}} \log \frac{\nu(x)}{\pi(x)} \nu(x) dx$ , where  $\nu, \pi \in \mathcal{P}(\mathcal{X})$  are absolutely continuous with respect to the Lebesgue measure on  $\mathcal{X}$  and  $\pi$  is a fixed reference probability measure on  $\mathcal{P}(\mathcal{X})$ . For other examples of functions h which satisfy the inequality  $TV^2(\nu',\nu) \leq \frac{1}{2}D_h(\nu',\nu)$ , and hence Assumption 2.1, see (Chizat, 2022b, Lemma 3.2) and (Kerimkulov et al., 2024, Proposition 2.18).

**Remark 2.3.** A similar notion to the Lipschitz property from Assumption 2.1, which goes under the name of Bregman continuity, was introduced in (Antonakopoulos et al., 2019) as a generalization of the standard Lipschitz continuity.

We are ready to state the first main result of the paper.

**Theorem 2.4** (Convergence of the simultaneous MDA algorithm (1)). Let  $(\nu^*, \mu^*)$  be an MNE of (1) and  $(\nu^0, \mu^0)$  be such that  $\sup_{\nu \in \mathcal{C}} D_h(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}} D_h(\mu, \mu^0) < \infty$ . Suppose that Assumption 1.1, 1.5, 1.6 and 2.1 hold. Suppose that  $\tau L \leq \frac{1}{2}$ , with  $L := \max\{L_{\nu}, L_{\mu}\}$ . Then, we have

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq 4\sqrt{\frac{L_{F}\left(\sup_{\nu\in\mathcal{C}}D_{h}(\nu,\nu^{0}) + \sup_{\mu\in\mathcal{D}}D_{h}(\mu,\mu^{0})\right)}{N}}.$$

**Remark 2.5.** Theorem 2.4 is consistent with the already known convergence rate  $\mathcal{O}\left(\frac{1}{\sqrt{N}}\right)$  of the MDA algorithm for min-max games with strategies in compact convex subsets of  $\mathbb{R}^d$ ; see e.g. (Bubeck, 2015, Theorem 5.1).

**Remark 2.6** (Initialization condition). The initialization requirement in Theorem 2.4, namely,  $\sup_{\nu \in \mathcal{C}} D_h(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}} D_h(\mu, \mu^0) < \infty$  must be verified case by case, depending on the choice of h and the admissible classes  $\mathcal{C}, \mathcal{D}$ . Such verifications for Examples 1.3 and 1.4 are carried out in Lemmas C.7 and C.8, respectively.

**Remark 2.7** (About the proof of Theorem 2.4). In their proof of convergence of the infinite-dimensional MD algorithm for convex F, (Aubin-Frankowski et al., 2022) show that relative smoothness is sufficient to prove that F is monotonically decreasing along the sequence  $(\nu^n)_{n\geq 0}$  generated by MD, i.e.,  $F(\nu^{n+1}) \leq F(\nu^n)$ , for all  $n \geq 0$ . The monotonicity property is key to establishing that the MD scheme converges to a minimizer of F with rate  $\mathcal{O}(\frac{1}{N})$ . In the case of (1), Assumption 1.6 and the fact that  $\tau L \leq \frac{1}{2}$  imply that  $F(\nu^{n+1}, \mu^n) \leq F(\nu^n, \mu^n) \leq F(\nu^n, \mu^{n+1})$ , for all  $n \geq 0$ . Thus, in the min-max setup, relative smoothness does not imply monotonic decay of F along the iterates. In contrast, we show that combining Assumption 1.6 with Assumption 2.1 allows us to control the Bregman divergence between consecutive iterates, i.e.,  $D_h(\nu^{n+1}, \nu^n)$  and  $D_h(\mu^{n+1}, \mu^n)$ , by  $\mathcal{O}(\tau^2)$  (see Lemma A.1). This condition will turn out to be sufficient to bypass the lack of monotonicity of F and also will guarantee the convergence in NI of the simultaneous MDA algorithm (1) with rate  $\mathcal{O}\left(\frac{1}{\sqrt{N}}\right)$ . For the proof, see Section A.

# 3 CONVERGENCE OF THE SEQUENTIAL MDA ALGORITHM (2)

Before we state the main result concerning the convergence of the sequential MDA algorithm (2), we introduce the necessary notions on the dual space of the space of probability measures.

Let  $(\mathcal{M}(\mathcal{X}), \|\cdot\|_{\mathrm{TV}})$  be the Banach space of finite signed measures  $\mu$  on  $\mathcal{X}$  equipped with the total variation norm  $\|\mu\|_{\mathrm{TV}} \coloneqq |\mu|(\mathcal{X})$ . Let  $(B_b(\mathcal{X}), \|\cdot\|_{\infty})$  be the Banach space of bounded measurable functions from  $\mathcal{X} \subset \mathbb{R}^d$  to  $(\mathbb{R}, |\cdot|)$ , where  $|\cdot|$  is the Euclidean norm. For any  $(f, m) \in B_b(\mathcal{X}) \times \mathcal{M}(\mathcal{X})$ , we define the duality pairing  $\langle \cdot, \cdot \rangle : B_b(\mathcal{X}) \times \mathcal{M}(\mathcal{X}) \to \mathbb{R}$  by

$$\langle f, m \rangle := \int_{\mathcal{X}} f(x) m(\mathrm{d}x).$$
 (5)

Next, we define the notion of convex conjugate of  $h: \mathcal{P}(\mathcal{X}) \to \mathbb{R}$  relative to the duality pairing (5). **Definition 3.1** (Convex conjugate). Let  $h: \mathcal{P}(\mathcal{X}) \to \mathbb{R}$  be a function. Then the map  $h^*: B_b(\mathcal{X}) \to \mathbb{R}$  given by  $h^*(f) := \sup_{m \in \mathcal{P}(\mathcal{X})} \{ \langle f, m \rangle - h(m) \}$  is called the convex conjugate of h.

Regardless of the convexity of h, it follows from (Bonnans & Shapiro, 2000, Theorem 2.112) that  $h^*$  is convex on  $B_b(\mathcal{X})$ , i.e., for all  $\lambda \in [0,1]$  and all  $f', f \in B_b(\mathcal{X})$ , we have that  $h^* ((1-\lambda)f + \lambda f') \leq (1-\lambda)h^*(f) + \lambda h^*(f')$ . In Example G.2, we provide the explicit form of  $h^*$  when h is the entropy. The following corollary shows that the first variation of  $h^*$  is the unique maximizer of  $m \mapsto \langle f, m \rangle - h(m)$ . This result is expected since on  $\mathbb{R}^d$  the "gradient" of the convex conjugate (of a strictly convex function) is the maximizer of the Legendre–Fenchel transformation.

**Corollary 3.2.** Let  $h^*: B_b(\mathcal{X}) \to \mathbb{R}$  be the convex conjugate of h. If Assumption 1.1 holds and  $h^*$  admits the first variation  $\frac{\delta h^*}{\delta f}(f)$  (cf. (42)) on  $B_b(\mathcal{X})$ , then

$$\frac{\delta h^*}{\delta f}(f) = \underset{m \in \mathcal{E}}{\arg \max} \left\{ \langle f, m \rangle - h(m) \right\}. \tag{6}$$

As shown in Example G.9, when h is chosen as the entropy, its convex conjugate  $h^*$  admits the first variation  $\frac{\delta h^*}{\delta f}(f)$ . If Assumption 1.1 holds, then, via (Hu et al., 2021, Lemma 4.1), we can characterize the convexity of  $h^*$  with respect to its first variation, i.e., for any  $f, f' \in B_b(\mathcal{X})$ ,

$$h^*(f') - h^*(f) \ge \int_{\mathcal{X}} (f'(x) - f(x)) \frac{\delta h^*}{\delta f}(f)(\mathrm{d}x),$$

and furthermore we can define the Bregman divergence between f and f' on the dual space.

**Definition 3.3** (Dual Bregman divergence). Let  $h^*: B_b(\mathcal{X}) \to \mathbb{R}$  be the convex conjugate of h. The dual  $h^*$ -Bregman divergence is the map  $D_{h^*}: B_b(\mathcal{X}) \times B_b(\mathcal{X}) \to [0, \infty)$  given by

$$D_{h^*}(f',f) := h^*(f') - h^*(f) - \int_{\mathcal{X}} (f'(x) - f(x)) \frac{\delta h^*}{\delta f}(f) (dx).$$

Since f,g are bounded and  $\frac{\delta h^*}{\delta g}(g)$  is a probability measure (cf. Definition G.8), it follows that  $\int_{\mathcal{X}} (f(x) - g(x)) \frac{\delta h^*}{\delta g}(g) (\mathrm{d}x)$  is well-defined. Moreover, since  $h^*$  is convex,  $D_{h^*}(f',f) \geq 0$ , for all  $f', f \in B_b(\mathcal{X})$ .

The following Lipschitzness assumption on the second variation  $\frac{\delta^2 h^*}{\delta f^2}$  (cf. Definition G.13) will turn out to be crucial for showing the improvement in the convergence rate of the sequential algorithm (2) compared to the simultaneous algorithm.

**Assumption 3.4.** Suppose that  $(B_b(\mathcal{X}) \times B_b(\mathcal{X})) \ni (f,g) \mapsto \frac{\delta^2 h^*}{\delta f^2}(f)(g) \in \mathcal{M}(\mathcal{X} \times \mathcal{X})$  is TV-Lipschitz, i.e., there exists  $L_{h^*} > 0$  such that, for all  $f, g, f', g' \in B_b(\mathcal{X})$ , it holds that

$$\left\| \frac{\delta^2 h^*}{\delta f^2} (f')(g') - \frac{\delta^2 h^*}{\delta f^2} (f)(g) \right\|_{\text{TV}} \le L_{h^*} \left( \|f' - f\|_{\infty} + \|g' - g\|_{\infty} \right),$$

In Example G.14 and Proposition G.15, we provide the explicit form of the second variation  $\frac{\delta^2 h^*}{\delta f^2}$  and verify Assumption 3.4, respectively, in the case where h is the entropy.

The following assumption ensures that F and its flat derivatives  $\frac{\delta F}{\delta \nu}$ ,  $\frac{\delta F}{\delta \mu}$  are uniformly bounded.

**Assumption 3.5** (Uniform boundedness of F and its flat derivatives). Suppose that there exists  $M>0, C_{\nu}>0$  and  $C_{\mu}>0$  such that for all  $\nu, \mu \in \mathcal{P}(\mathcal{X})$ , and all  $x,y \in \mathcal{X}$ , we have

$$|F(\nu,\mu)| \le M, \quad \left| \frac{\delta F}{\delta \nu}(\nu,\mu,x) \right| \le C_{\nu}, \quad \left| \frac{\delta F}{\delta \mu}(\nu,\mu,y) \right| \le C_{\mu}.$$

In Proposition D.1 and E.1, we verify that Assumption 3.5 is satisfied by Example 1.2 and E.

Now, we are ready to state the second main result of the paper.

**Theorem 3.6** (Convergence of the sequential MDA algorithm (2)). Let  $(\nu^*, \mu^*)$  be an MNE of (1) and  $(\nu^0, \mu^0)$  be such that  $\sup_{\nu \in \mathcal{C}} D_h(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}} D_h(\mu, \mu^0) < \infty$  (cf. Remark 2.6). Let Assumption 1.1, 1.5, 1.6, 2.1, 3.4 and 3.5 hold. Suppose that  $\tau L \leq \frac{1}{2}$ , with  $L := \max\{L_{\nu}, L_{\mu}\}$ . Then, we have

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq \frac{1}{2N^{2/3}} \left(3\left(\sup_{\nu\in\mathcal{C}}D_{h}(\nu,\nu^{0}) + \sup_{\mu\in\mathcal{D}}D_{h}(\mu,\mu^{0})\right)^{2/3} \times \left(\kappa L_{h^{*}} + 2L_{F}L\right)\right)^{1/3} + 2M\right), \quad (7)$$

where  $\kappa := C_{\nu}^3 + C_{\mu}^3$ .

**Remark 3.7.** In particular, if  $F(\nu,\mu) = \int_{\mathcal{X}} \int_{\mathcal{X}} f(x,y)\nu(\mathrm{d}x)\mu(\mathrm{d}y)$ , for a bounded function  $f:\mathcal{X}\times\mathcal{X}\to\mathbb{R}$ , then Assumption 2.1 is satisfied and in Definition 1.6 we have  $L_{\nu}=L_{\mu}=0$ . Therefore, L=0 in (7), and hence Theorem 3.6 is consistent with the already known convergence rate  $\mathcal{O}\left(\frac{1}{N^{2/3}}\right)$  of the MDA algorithm for min-max games with strategies in compact convex subsets of  $\mathbb{R}^d$  and bilinear payoff function; see (Wibisono et al., 2022, Theorem 3.2 and Corollary 3.3). Since we work in an infinite-dimensional setting with a non-linear convex-concave objective function F, Theorem 3.6 substantially generalizes the results of (Wibisono et al., 2022).

Remark 3.8 (About the proof of Theorem 3.6). The main difference compared to Theorem 2.4 is the extra term  $F(\nu^{n+1},\mu^n)-F(\nu^n,\mu^n)$  which arises from the non-symmetry of the flat derivatives of F in Algorithm (2). We combine this difference with  $\int_{\mathcal{X}} \frac{\delta F}{\delta \nu}(\nu^n,\mu^n,x)(\nu^n-\nu^{n+1})(\mathrm{d}x)$  and  $\int_{\mathcal{X}} \frac{\delta F}{\delta \mu}(\nu^{n+1},\mu^n,y)(\mu^{n+1}-\mu^n)(\mathrm{d}y)$  via relative smoothness. This produces the Bregman commutators  $D_h(\nu^n,\nu^{n+1})-D_h(\nu^{n+1},\nu^n)$  and  $D_h(\mu^n,\mu^{n+1})-D_h(\mu^{n+1},\mu^n)$ . To handle these commutators, we pass from the measure space to the dual space of bounded measurable functions. Hence, the commutators become  $D_{h^*}\left(\frac{\delta h}{\delta \nu}(\nu^{n+1},\cdot),\frac{\delta h}{\delta \nu}(\nu^n,\cdot)\right)-D_{h^*}\left(\frac{\delta h}{\delta \nu}(\nu^n,\cdot),\frac{\delta h}{\delta \nu}(\nu^{n+1},\cdot)\right)$  and analogously for  $\mu$ . Applying Assumption 3.4 and 3.5, we show that these commutators are of order  $\mathcal{O}(\tau^3)$ . This refined estimate yields the improved convergence rate  $\mathcal{O}\left(\frac{1}{N^{2/3}}\right)$ . For the proof, see Section A.

# 4 NUMERICAL EXAMPLE

In this section, we outline how to implement the infinite-dimensional algorithms (1) and (2) in the case where h is the relative entropy. For brevity, we present the derivations only for Algorithm (1), as the arguments for Algorithm (2) are entirely analogous. The complete algorithms for both the simultaneous and sequential MDA schemes can be found in Algorithm (3) and Algorithm (4) in Section B.

#### 4.1 SIMULATION OF INFINITE-DIMENSIONAL MDA

As shown in Example 1.3, by taking h to be the entropy, the corresponding h-Bregman divergence is exactly the KL divergence. Moreover, using the flat derivative formula (4), the first-order optimality condition (Hu et al., 2021, Proposition 2.5) applied to  $(\nu^{n+1}, \mu^{n+1})$  in Algorithm (1) gives

$$\begin{cases} \log \nu^{n+1}(x) - \log \nu^n(x) = -\tau \frac{\delta F}{\delta \nu}(\nu^n, \mu^n, x) + C, \\ \log \mu^{n+1}(y) - \log \mu^n(y) = \tau \frac{\delta F}{\delta \mu}(\nu^n, \mu^n, y) + C', \end{cases}$$

for every  $n \geq 0$  and, for all  $x, y \in \mathcal{X}$  Lebesgue a.e., where  $C, C' \in \mathbb{R}$ . By summing over n and exponentiating both sides, we obtain

$$\begin{cases} \nu^{n}(x) \propto \nu^{0}(x)e^{-\tau \sum_{k=0}^{n-1} \frac{\delta F}{\delta \nu}} (\nu^{k}, \mu^{k}, x), \\ \mu^{n}(y) \propto \mu^{0}(y)e^{\tau \sum_{k=0}^{n-1} \frac{\delta F}{\delta \mu}} (\nu^{k}, \mu^{k}, y), \end{cases}$$

where the constants C,C' are absorbed into the normalizations. For simplicity, suppose the initial samples  $(X_j,Y_j)_{j=1}^J$  are drawn uniformly, so that  $(\nu^0,\mu^0)$  are uniform densities. We set  $(X_{j,0},Y_{j,0})_{j=1}^J=(X_j,Y_j)_{j=1}^J$  and sample from  $(\nu^1,\mu^1)$  via Langevin dynamics:

$$X_{j,t+1} = X_{j,t} - \gamma \nabla \frac{\delta F}{\delta \nu}(\nu^0, \mu^0, X_{j,t}) + \sqrt{\frac{2\gamma}{\tau}} \mathcal{N}_{j,t},$$

$$Y_{j,t+1} = Y_{j,t} + \gamma \nabla \frac{\delta F}{\delta \mu}(\nu^0, \mu^0, Y_{j,t}) + \sqrt{\frac{2\gamma}{\tau}} \mathcal{N}_{j,t},$$

for  $1 \leq j \leq J$  and  $0 \leq t \leq T-1$ , where  $\gamma > 0$  is the step size and  $\mathcal{N}_{j,t}$  are i.i.d standard Gaussian variables. For sufficiently large J and T, the terminal particles  $(X_{j,T},Y_{j,T})_{j=1}^J$  approximate samples from  $(\nu^1,\mu^1)$ . Repeating this procedure recursively then yields samples from  $(\nu^2,\mu^2),\ldots,(\nu^n,\mu^n)$ .

# 4.2 TRAINING GANS BY MDA

We train the mean-field GAN from Example 1.2 using simultaneous and sequential MDA-GAN (Algorithms (5) and (6)) on the 8-Gaussian mixture and Swiss Roll datasets (Gulrajani et al., 2017). Full algorithmic details, including hyperparameters and network architectures, are in Section B. Both methods are run for 2000 iterations, with performance assessed by visualizing generated samples at 400, 1000, and 2000 iterations.

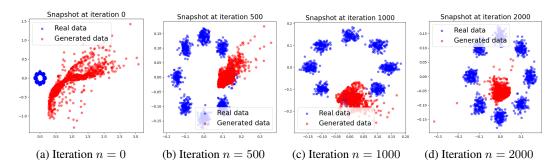


Figure 1: Simultaneous MDA-GAN (Algorithm 5) learning an 8-Gaussian mixture

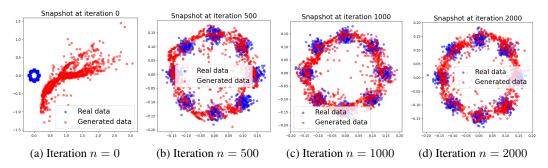


Figure 2: Sequential MDA-GAN (Algorithm 6) learning an 8-Gaussian mixture

Figures 1 and 2 show the training dynamics of simultaneous and sequential MDA-GANs on the 8-Gaussian mixture, with analogous results on the Swiss Roll in Figures 3 and 4. In both settings, generated samples start far from the data but the sequential variant captures the multi-modal structure and the spiral geometry of the Swiss Roll more clearly and at earlier iterations. In Section B, we plot the  $L^1$ -Wasserstein distance  $W_1(T_{\theta^n}\#\xi,\hat{\xi})$  for both tasks over iterations n, confirming the faster convergence of sequential MDA-GAN.

# REFERENCES

- R. Abraham, J.E. Marsden, and T. Ratiu. *Manifolds, Tensor Analysis, and Applications*. Applied Mathematical Sciences. Springer New York, 2012.
  - C.D. Aliprantis and K.C. Border. *Infinite Dimensional Analysis: A Hitchhiker's Guide*. Springer, 2007.
- A. Ambrosetti and G. Prodi. *A Primer of Nonlinear Analysis*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 1995.
  - Luigi Ambrosio, Nicola Fusco, and Diego Pallara. Functions of Bounded Variation and Free Discontinuity Problems. Oxford University Press, 2000.
  - Mihai Anitescu. Degenerate nonlinear programming with a quadratic growth condition. *SIAM Journal on Optimization*, 10(4):1116–1135, 2000.
  - Kimon Antonakopoulos, Veronica Belmega, and Panayotis Mertikopoulos. An adaptive mirror-prox method for variational inequalities with singular operators. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
  - Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 214–223. PMLR, 06–11 Aug 2017.
  - Pierre-Cyril Aubin-Frankowski, Anna Korba, and Flavien Léger. Mirror descent with relative smoothness in measure spaces, with application to Sinkhorn and EM. In *Advances in Neural Information Processing Systems*, volume 35, pp. 17263–17275. Curran Associates, Inc., 2022.
  - Xingjian Bai, Guangyi He, Yifan Jiang, and Jan Obloj. Wasserstein distributional robustness of neural networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
  - Heinz H. Bauschke, Jérôme Bolte, and Marc Teboulle. A descent lemma beyond Lipschitz gradient continuity: First-order methods revisited and applications. *Math. Oper. Res.*, 42:330–348, 2017.
  - Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
  - J. Frédéric Bonnans and Alexander Shapiro. Perturbation Analysis of Optimization Problems. Springer Series in Operations Research, 2000.
  - Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Found. Trends Mach. Learn.*, 8(3–4):231–357, 2015.
  - René A. Carmona and François Delarue. *Probabilistic Theory of Mean Field Games with Applications I: Mean Field FBSDEs, Control, and Games.* Springer International Publishing, 2018.
  - Lénaïc Chizat. Mean-field langevin dynamics: Exponential convergence and annealing. *Transactions on Machine Learning Research*, 2022a.
  - Lénaïc Chizat. Convergence Rates of Gradient Methods for Convex Optimization in the Space of Measures. *Open Journal of Mathematical Optimization*, 3:8, 2022b.
  - Lénaïc Chizat and Francis R. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In *NeurIPS*, 2018.
- Philippe G. Ciarlet. *Linear and Nonlinear Functional Analysis with Applications*. Society for Industrial and Applied Mathematics, 2013.
- Carles Domingo-Enrich, Samy Jelassi, Arthur Mensch, Grant Rotskoff, and Joan Bruna. A mean-field analysis of two-player zero-sum games. In *Advances in Neural Information Processing Systems*, volume 33, pp. 20215–20226. Curran Associates, Inc., 2020.
  - P. Dupuis and R. S. Ellis. A Weak Convergence Approach to the Theory of Large Deviations. Wiley, New York, NY, 1997.

- Pavel Dvurechensky and Jia-Jie Zhu. Analysis of kernel mirror prox for measure optimization. In 27th International Conference on Artificial Intelligence and Statistics (AISTATS), 2024.
  - I. L. Glicksberg. A further generalization of the kakutani fixed point theorem, with application to nash equilibrium points. *Proceedings of the American Mathematical Society*, 3(1):170–174, 1952.
  - Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved training of Wasserstein GANs, 2017. arXiv:1704.00028.
  - Ya-Ping Hsieh, Chen Liu, and Volkan Cevher. Finding mixed Nash equilibria of generative adversarial networks. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2810–2819. PMLR, 09–15 Jun 2019.
  - Kaitong Hu, Zhenjie Ren, David Šiška, and Łukasz Szpruch. Mean-field Langevin dynamics and energy landscape of neural networks. *Annales de l'Institut Henri Poincaré*, *Probabilités et Statistiques*, 57(4):2043 2065, 2021.
  - Bekzhan Kerimkulov, James-Michael Leahy, David Šiška, Łukasz Szpruch, and Yufei Zhang. A Fisher-Rao gradient flow for entropy-regularised Markov decision processes in Polish spaces, 2023. arXiv:2310.02951.
  - Bekzhan Kerimkulov, David Šiška, Łukasz Szpruch, and Yufei Zhang. Mirror Descent for Stochastic Control Problems with Measure-valued Controls, 2024. arXiv:2401.01198.
  - Juno Kim, Kakei Yamamoto, Kazusato Oko, Zhuoran Yang, and Taiji Suzuki. Symmetric mean-field Langevin dynamics for distributional minimax problems. In *The Twelfth International Conference on Learning Representations*, 2024.
  - Razvan-Andrei Lascu, Mateusz Majka, and Łukasz Szpruch. A Fisher-Rao gradient flow for entropic mean-field min-max games. *Transactions on Machine Learning Research*, 2024.
  - Razvan-Andrei Lascu, Mateusz B. Majka, and Łukasz Szpruch. Entropic mean-field min-max problems via Best Response flow. *Appl. Math. Optim.*, 91(2), March 2025.
  - Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, and Barnabas Poczos. MMD GAN: Towards deeper understanding of moment matching network. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
  - Linshan Liu, Mateusz B. Majka, and Łukasz Szpruch. Polyak–Łojasiewicz inequality on the space of measures and convergence of mean-field birth-death processes. *Applied Mathematics and Optimization*, 87(3), 2023.
  - Haihao Lu, Robert M. Freund, and Yurii Nesterov. Relatively smooth convex optimization by first-order methods, and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.
  - Yulong Lu. Two-scale gradient descent ascent dynamics finds mixed Nash equilibria of continuous games: a mean-field perspective. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
  - Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences of the United States of America*, 115:E7665 E7671, 2018.
  - A.S. Nemirovski and D.B. Yudin. *Problem Complexity and Method Efficiency in Optimization*. A Wiley-Interscience publication. Wiley, 1983.
  - Hukukane Nikaidô and Kazuo Isoda. Note on non-cooperative convex game. *Pacific Journal of Mathematics*, 5:807–815, 1955.
  - Atsushi Nitanda, Denny Wu, and Taiji Suzuki. Convex analysis of the mean field Langevin dynamics. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, 2022.

- J.M. Ortega and W.C. Rheinboldt. *Iterative Solution of Nonlinear Equations in Several Variables*. Classics in Applied Mathematics. Society for Industrial and Applied Mathematics, 1970.
- Grant M. Rotskoff and Eric Vanden-Eijnden. Neural networks as interacting particle systems: Asymptotic convexity of the loss landscape and universal scaling of the approximation error, 2018. arXiv:1805.00915.
- F. Santambrogio. *Optimal Transport for Applied Mathematicians: Calculus of Variations, PDEs, and Modeling.* Progress in Nonlinear Differential Equations and Their Applications. Springer International Publishing, 2015. ISBN 9783319208282.
- Nian Si, Fan Zhang, Zhengyuan Zhou, and Jose Blanchet. Distributionally robust batch contextual bandits. *Management Science*, 69(10):5772–5793, 2023.
- Maurice Sion. On general minimax theorems. *Pacific Journal of Mathematics*, 8(1):171 176, 1958.
- Taiji Suzuki, Denny Wu, and Atsushi Nitanda. Mean-field Langevin dynamics: Time-space discretization, stochastic gradient, and variance reduction. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Manan Tomar, Lior Shani, Yonathan Efroni, and Mohammad Ghavamzadeh. Mirror Descent Policy Optimization, 2021. arXiv:2005.09814.
- Camilo Garcia Trillos and Nicolas Garcia Trillos. On adversarial robustness and the use of Wasserstein ascent-descent dynamics to enforce it, 2023. arXiv:2301.03662.
- A.B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Series in Statistics. Springer New York, 2008.
- Guillaume Wang and Lénaïc Chizat. An exponentially converging particle method for the mixed Nash equilibrium of continuous games, 2023. arXiv:2211.01280.
- Andre Wibisono, Molei Tao, and Georgios Piliouras. Alternating Mirror Descent for constrained min-max games. In *Advances in Neural Information Processing Systems*, volume 35, pp. 35201–35212. Curran Associates, Inc., 2022.
- C. Zalinescu. Convex Analysis In General Vector Spaces. World Scientific Publishing Company, 2002.

# A PROOFS OF THEOREM 2.4 AND THEOREM 3.6

- This section is dedicated to the proofs of the main results, namely Theorem 2.4 and Theorem 3.6. Before we proceed, we will need an auxiliary result, which will turn out to be essential for proving both main theorems. The proof of Lemma A.1 is given in Appendix C.
- **Lemma A.1.** Let Assumption 1.1, 1.6 and 2.1 hold. Suppose that  $\tau L \leq \frac{1}{2}$ , with  $L := \max\{L_{\nu}, L_{\mu}\}$ . Then, for both Algorithms (1) and (2), it holds, for all  $n \geq 0$ , that

$$D_h(\nu^{n+1}, \nu^n) \le 4L_F \tau^2$$
 and  $D_h(\mu^{n+1}, \mu^n) \le 4L_F \tau^2$ .

#### A.1 Proof of Theorem 2.4

*Proof of Theorem 2.4.* Since  $\nu \mapsto \tau \int_{\mathcal{X}} \frac{\delta F}{\delta \nu}(\nu^n, \mu^n, x)(\nu - \nu^n)(\mathrm{d}x)$  is convex, applying Lemma C.1 with  $\bar{\nu} = \nu^{n+1}$  and  $\mu = \nu^n$  implies that, for any  $\nu \in \mathcal{C}$ , we have

$$\tau \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) (\nu - \nu^{n}) (\mathrm{d}x) + D_{h}(\nu, \nu^{n}) \ge \tau \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) (\nu^{n+1} - \nu^{n}) (\mathrm{d}x) + D_{h}(\nu^{n+1}, \nu^{n}) + D_{h}(\nu, \nu^{n+1}),$$

or, equivalently,

$$-\tau \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) (\nu - \nu^{n}) (\mathrm{d}x) - D_{h}(\nu, \nu^{n}) \leq -\tau \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) (\nu^{n+1} - \nu^{n}) (\mathrm{d}x) - D_{h}(\nu^{n+1}, \nu^{n}) - D_{h}(\nu, \nu^{n+1}).$$
 (8)

Similarly, since  $\mu \mapsto -\tau \int_{\mathcal{X}} \frac{\delta F}{\delta \mu}(\nu^n, \mu^n, y)(\mu - \mu^n)(\mathrm{d}y)$  is convex, applying Lemma C.1 with  $\bar{\nu} = \mu^{n+1}$  and  $\mu = \mu^n$  implies that, for any  $\mu \in \mathcal{D}$ , we have

$$\tau \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n}, \mu^{n}, y) (\mu - \mu^{n}) (\mathrm{d}y) - D_{h}(\mu, \mu^{n}) \le \tau \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n}, \mu^{n}, y) (\mu^{n+1} - \mu^{n}) (\mathrm{d}y) - D_{h}(\mu^{n+1}, \mu^{n}) - D_{h}(\mu, \mu^{n+1}). \tag{9}$$

Using the convexity of  $\nu \mapsto F(\nu, \mu)$  in (8), with  $\nu = \nu^n$  and  $\mu = \mu^n$ , we have that

$$F(\nu^{n}, \mu^{n}) - F(\nu, \mu^{n}) - \frac{1}{\tau} D_{h}(\nu, \nu^{n}) \leq \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) (\nu^{n} - \nu^{n+1}) (\mathrm{d}x) - \frac{1}{\tau} D_{h}(\nu^{n+1}, \nu^{n}) - \frac{1}{\tau} D_{h}(\nu, \nu^{n+1}).$$
 (10)

From  $L_{\nu}$ -relative smoothness and the fact that  $\tau L \leq \frac{1}{2} < 1$ , it follows that

$$F(\nu^{n+1}, \mu^n) \le F(\nu^n, \mu^n) + \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^n, \mu^n, x) (\nu^{n+1} - \nu^n) (\mathrm{d}x) + L_{\nu} D_h(\nu^{n+1}, \nu^n)$$

$$\le F(\nu^n, \mu^n) + \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^n, \mu^n, x) (\nu^{n+1} - \nu^n) (\mathrm{d}x) + \frac{1}{\tau} D_h(\nu^{n+1}, \nu^n).$$
 (11)

Hence, combining (10) with (11), we obtain for any  $\nu \in \mathcal{C}$  that

$$F(\nu^n, \mu^n) - F(\nu, \mu^n) - \frac{1}{\tau} D_h(\nu, \nu^n) \le F(\nu^n, \mu^n) - F(\nu^{n+1}, \mu^n) - \frac{1}{\tau} D_h(\nu, \nu^{n+1}). \tag{12}$$

Similarly, using concavity of  $\mu \mapsto F(\nu, \mu)$  in (9), with  $\nu = \nu^n$  and  $\mu = \mu^n$ , we have that

$$F(\nu^{n}, \mu) - F(\nu^{n}, \mu^{n}) - \frac{1}{\tau} D_{h}(\mu, \mu^{n}) \leq \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n}, \mu^{n}, y) (\mu^{n+1} - \mu^{n}) (\mathrm{d}y) - \frac{1}{\tau} D_{h}(\mu^{n+1}, \mu^{n}) - \frac{1}{\tau} D_{h}(\mu, \mu^{n+1}).$$
(13)

From  $L_{\mu}$ -relative smoothness and the fact that  $\tau L \leq \frac{1}{2} < 1$ , it follows that

$$F(\nu^{n}, \mu^{n+1}) \geq F(\nu^{n}, \mu^{n}) + \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n}, \mu^{n}, y) (\mu^{n+1} - \mu^{n}) (\mathrm{d}y) - L_{\mu} D_{h}(\mu^{n+1}, \mu^{n})$$

$$\geq F(\nu^{n}, \mu^{n}) + \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n}, \mu^{n}, y) (\mu^{n+1} - \mu^{n}) (\mathrm{d}y) - \frac{1}{\tau} D_{h}(\mu^{n+1}, \mu^{n}). \quad (14)$$

Hence, combining (13) with (14), we obtain for any  $\mu \in \mathcal{D}$  that

$$F(\nu^{n},\mu) - F(\nu^{n},\mu^{n}) - \frac{1}{\tau}D_{h}(\mu,\mu^{n}) \leq F(\nu^{n},\mu^{n+1}) - F(\nu^{n},\mu^{n}) - \frac{1}{\tau}D_{h}(\mu,\mu^{n+1}).$$
 (15)

Adding inequalities (12) and (15) implies that for any  $(\nu, \mu) \in \mathcal{C} \times \mathcal{D}$  we have

$$F(\nu^{n},\mu) - F(\nu,\mu^{n}) \leq F(\nu^{n},\mu^{n}) - F(\nu^{n+1},\mu^{n}) + F(\nu^{n},\mu^{n+1}) - F(\nu^{n},\mu^{n}) + \frac{1}{\tau}D_{h}(\nu,\nu^{n}) + \frac{1}{\tau}D_{h}(\mu,\mu^{n}) - \frac{1}{\tau}D_{h}(\nu,\nu^{n+1}) - \frac{1}{\tau}D_{h}(\mu,\mu^{n+1}).$$
 (16)

By Assumption 2.1, we have that

$$|F(\nu^n,\mu^n) - F(\nu^{n+1},\mu^n)|^2 = |F(\nu^{n+1},\mu^n) - F(\nu^n,\mu^n)|^2 \le L_F D_h(\nu^{n+1},\nu^n) \le 4L_F^2 \tau^2,$$

and

$$|F(\nu^n, \mu^{n+1}) - F(\nu^n, \mu^n)|^2 \le L_F D_h(\mu^{n+1}, \mu^n) \le 4L_F^2 \tau^2,$$

where the last inequalities follows from Lemma A.1. Therefore, from (16), we obtain

$$F(\nu^n, \mu) - F(\nu, \mu^n) \le 4L_F \tau + \frac{1}{\tau} D_h(\nu, \nu^n) + \frac{1}{\tau} D_h(\mu, \mu^n) - \frac{1}{\tau} D_h(\nu, \nu^{n+1}) - \frac{1}{\tau} D_h(\mu, \mu^{n+1}),$$

Summing the previous inequality over n=0,1,...,N-1, using  $D_h(\nu,\nu^N)+D_h(\mu,\mu^N)\geq 0$ , for any  $(\nu,\mu)\in\mathcal{C}\times\mathcal{D}$ , bounding the right-hand from above by its supremum over  $(\nu,\mu)\in\mathcal{C}\times\mathcal{D}$ , and dividing by N gives

$$\frac{1}{N} \sum_{n=0}^{N-1} \left( F(\nu^n, \mu) - F(\nu, \mu^n) \right) \le 4L_F \tau + \frac{1}{N\tau} \left( \sup_{\nu \in \mathcal{C}} D_h(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}} D_h(\mu, \mu^0) \right). \tag{17}$$

Since  $\nu \mapsto F(\nu,\mu)$  and  $\mu \mapsto -F(\nu,\mu)$  are convex, it follows by Jensen's inequality that

$$\frac{1}{N} \sum_{n=0}^{N-1} \left( F(\nu^n, \mu) - F(\nu, \mu^n) \right) = \frac{1}{N} \sum_{n=0}^{N-1} F(\nu^n, \mu) - \frac{1}{N} \sum_{n=0}^{N-1} F(\nu, \mu^n) \\
\ge F\left(\frac{1}{N} \sum_{n=0}^{N-1} \nu^n, \mu\right) - F\left(\nu, \frac{1}{N} \sum_{n=0}^{N-1} \mu^n\right). \tag{18}$$

Combining (17) with (18) and taking maximum over  $(\nu, \mu)$  gives

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq 4L_{F}\tau + \frac{1}{N\tau}\left(\sup_{\nu \in \mathcal{C}}D_{h}(\nu, \nu^{0}) + \sup_{\mu \in \mathcal{D}}D_{h}(\mu, \mu^{0})\right).$$

Minimizing the right-hand side over  $\tau$  amounts to taking

$$\tau = \frac{1}{2} \sqrt{\frac{\sup_{\nu \in \mathcal{C}} D_h(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}} D_h(\mu, \mu^0)}{L_F N}},$$

and hence we obtain

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq 4\sqrt{\frac{L_{F}\left(\sup_{\nu\in\mathcal{C}}D_{h}(\nu,\nu^{0}) + \sup_{\mu\in\mathcal{D}}D_{h}(\mu,\mu^{0})\right)}{N}}.$$

# A.2 Proof of Theorem 3.6

*Proof of Theorem 3.6.* We start the proof by following the same calculations from Theorem 2.4. For (2), after applying Lemma C.1 and using convexity-concavity of F, (10) remains unchanged, i.e.,

$$F(\nu^{n}, \mu^{n}) - F(\nu, \mu^{n}) - \frac{1}{\tau} D_{h}(\nu, \nu^{n}) \leq \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) (\nu^{n} - \nu^{n+1}) (\mathrm{d}x) - \frac{1}{\tau} D_{h}(\nu^{n+1}, \nu^{n}) - \frac{1}{\tau} D_{h}(\nu, \nu^{n+1}),$$

while (13) becomes

$$F(\nu^{n+1}, \mu) - F(\nu^{n+1}, \mu^n) - \frac{1}{\tau} D_h(\mu, \mu^n) \le \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n+1}, \mu^n, y) (\mu^{n+1} - \mu^n) (\mathrm{d}y) - \frac{1}{\tau} D_h(\mu^{n+1}, \mu^n) - \frac{1}{\tau} D_h(\mu, \mu^{n+1}).$$

Adding the previous two inequalities, summing the resulting inequality over n=0,1,...,N-1, dividing by N, using (18) and taking maximum over  $(\nu,\mu)$  we arrive at

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq \frac{1}{N}\sum_{n=0}^{N-1}\left(\int_{\mathcal{X}}\frac{\delta F}{\delta\nu}(\nu^{n}, \mu^{n}, x)(\nu^{n} - \nu^{n+1})(\mathrm{d}x)\right) + \int_{\mathcal{X}}\frac{\delta F}{\delta\mu}(\nu^{n+1}, \mu^{n}, y)(\mu^{n+1} - \mu^{n})(\mathrm{d}y) + \frac{1}{N\tau}\left(\sup_{\nu \in \mathcal{C}}D_{h}(\nu, \nu^{0}) + \sup_{\mu \in \mathcal{D}}D_{h}(\mu, \mu^{0})\right) + \frac{1}{N}\sum_{n=0}^{N-1}\left(F(\nu^{n+1}, \mu^{n}) - F(\nu^{n}, \mu^{n})\right) - \frac{1}{N\tau}\sum_{n=0}^{N-1}\left(D_{h}(\nu^{n+1}, \nu^{n}) + D_{h}(\mu^{n+1}, \mu^{n})\right), \quad (19)$$

where we used the fact that  $D_h(\nu, \nu^N) + D_h(\mu, \mu^N) \ge 0$ , for any  $(\nu, \mu) \in \mathcal{C} \times \mathcal{D}$ .

Note that the key difference to the estimates from Theorem 2.4 is the appearance of the term  $F(\nu^{n+1},\mu^n)-F(\nu^n,\mu^n)$  due to the non-symmetry of the flat derivatives of F in (2). The idea is to combine  $F(\nu^{n+1},\mu^n)-F(\nu^n,\mu^n)$  with both  $\int_{\mathcal{X}} \frac{\delta F}{\delta \nu}(\nu^n,\mu^n,x)(\nu^n-\nu^{n+1})(\mathrm{d}x)$  and  $\int_{\mathcal{X}} \frac{\delta F}{\delta \mu}(\nu^{n+1},\mu^n,y)(\mu^{n+1}-\mu^n)(\mathrm{d}y)$  via relative smoothness in order to obtain  $D_h(\nu^n,\nu^{n+1})-D_h(\nu^{n+1},\nu^n)$  and  $D_h(\mu^n,\mu^{n+1})-D_h(\mu^{n+1},\mu^n)$ , which will prove to be of order  $\mathcal{O}(\tau^3)$ .

Since the flat derivative of  $\mathcal{E} \ni m \mapsto D_h(m,m') \in [0,\infty)$  is given by  $\frac{\delta}{\delta m} D_h(\cdot,m') = \frac{\delta h}{\delta m}(m,x) - \frac{\delta h}{\delta m}(m',x)$ , it follows that the first-order conditions for (2) read

$$\begin{cases} \frac{\delta h}{\delta \nu}(\nu^{n+1}, x) - \frac{\delta h}{\delta \nu}(\nu^{n}, x) = -\tau \frac{\delta F}{\delta \nu}(\nu^{n}, \mu^{n}, x), \\ \frac{\delta h}{\delta \mu}(\mu^{n+1}, y) - \frac{\delta h}{\delta \mu}(\mu^{n}, y) = \tau \frac{\delta F}{\delta \mu}(\nu^{n+1}, \mu^{n}, y), \end{cases}$$
(20)

for all  $(x,y) \in \mathcal{X} \times \mathcal{X}$  Lebesgue a.e. It can be shown directly from Definition 1.2 that

$$\int_{\mathcal{X}} \left( \frac{\delta h}{\delta \nu} (\nu', x) - \frac{\delta h}{\delta \nu} (\nu, x) \right) (\nu' - \nu) (\mathrm{d}x) = D_h(\nu', \nu) + D_h(\nu, \nu'), \tag{21}$$

for all  $\nu, \nu' \in \mathcal{C}$ , and analogously for  $D_h(\mu', \mu) + D_h(\mu, \mu')$ . Then, using (20) and (21) we obtain that

$$-\int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^n, \mu^n, x) (\nu^{n+1} - \nu^n) (\mathrm{d}x) = \frac{1}{\tau} \int_{\mathcal{X}} \left( \frac{\delta h}{\delta \nu} (\nu^{n+1}, x) - \frac{\delta h}{\delta \nu} (\nu^n, x) \right) (\nu^{n+1} - \nu^n) (\mathrm{d}x)$$
$$= \frac{1}{\tau} \left( D_h(\nu^{n+1}, \nu^n) + D_h(\nu^n, \nu^{n+1}) \right), \quad (22)$$

and similarly

$$\int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n+1}, \mu^{n}, y) (\mu^{n+1} - \mu^{n}) (\mathrm{d}y) = \frac{1}{\tau} \int_{\mathcal{X}} \left( \frac{\delta h}{\delta \mu} (\mu^{n+1}, y) - \frac{\delta h}{\delta \mu} (\mu^{n}, y) \right) (\mu^{n+1} - \mu^{n}) (\mathrm{d}y) \\
= \frac{1}{\tau} \left( D_{h} (\mu^{n+1}, \mu^{n}) + D_{h} (\mu^{n}, \mu^{n+1}) \right). \tag{23}$$

Therefore, using (22) and (23) in (19), we obtain that

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq \frac{1}{N\tau}\left(\sup_{\nu\in\mathcal{C}}D_{h}(\nu,\nu^{0}) + \sup_{\mu\in\mathcal{D}}D_{h}(\mu,\mu^{0})\right) + \frac{1}{N\tau}\sum_{n=0}^{N-1}\left(D_{h}(\nu^{n+1},\nu^{n}) + D_{h}(\nu^{n},\nu^{n+1}) + D_{h}(\mu^{n+1},\mu^{n}) + D_{h}(\mu^{n},\mu^{n+1})\right) + \frac{1}{N}\sum_{n=0}^{N-1}\left(F(\nu^{n+1},\mu^{n}) - F(\nu^{n},\mu^{n})\right) - \frac{1}{N\tau}\sum_{n=0}^{N-1}\left(D_{h}(\nu^{n+1},\nu^{n}) + D_{h}(\mu^{n+1},\mu^{n})\right).$$
(24)

Then, we observe that

$$D_h(\nu^n, \nu^{n+1}) = \frac{1}{2} \left( D_h(\nu^n, \nu^{n+1}) - D_h(\nu^{n+1}, \nu^n) \right) + \frac{1}{2} \left( D_h(\nu^n, \nu^{n+1}) + D_h(\nu^{n+1}, \nu^n) \right), \tag{25}$$

and a similar representation holds for  $D_h(\mu^n, \mu^{n+1})$ . Similarly, we can write

$$F(\nu^{n+1}, \mu^n) - F(\nu^n, \mu^n) = \frac{1}{2} \left( F(\nu^{n+1}, \mu^n) - F(\nu^n, \mu^n) \right)$$

$$+ \frac{1}{2} \left( F(\nu^{n+1}, \mu^n) - F(\nu^{n+1}, \mu^{n+1}) + F(\nu^{n+1}, \mu^{n+1}) - F(\nu^n, \mu^n) \right). \tag{26}$$

Therefore, putting (25) and (26) into (24) gives

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq \frac{1}{N\tau}\left(\sup_{\nu\in\mathcal{C}}D_{h}(\nu,\nu^{0}) + \sup_{\mu\in\mathcal{D}}D_{h}(\mu,\mu^{0})\right) + \frac{1}{2N\tau}\sum_{n=0}^{N-1}\left(D_{h}(\nu^{n},\nu^{n+1}) - D_{h}(\nu^{n+1},\nu^{n}) + D_{h}(\mu^{n},\mu^{n+1}) - D_{h}(\mu^{n+1},\mu^{n})\right) + \frac{1}{2N}\sum_{n=0}^{N-1}\left(\frac{1}{\tau}\left(D_{h}(\nu^{n},\nu^{n+1}) + D_{h}(\nu^{n+1},\nu^{n})\right) + F(\nu^{n+1},\mu^{n}) - F(\nu^{n},\mu^{n})\right) + \frac{1}{2N}\sum_{n=0}^{N-1}\left(\frac{1}{\tau}\left(D_{h}(\mu^{n},\mu^{n+1}) + D_{h}(\mu^{n+1},\mu^{n})\right) + F(\nu^{n+1},\mu^{n}) - F(\nu^{n+1},\mu^{n+1}) + F(\nu^{n+1},\mu^{n+1}) - F(\nu^{n},\mu^{n})\right).$$

$$\left(27\right)$$

Combining the fact that  $\nu \mapsto F(\nu, \mu)$  is  $L_{\nu}$ -smooth relative to h with the first-order condition (20), we have that

$$F(\nu^{n+1}, \mu^{n}) - F(\nu^{n}, \mu^{n}) \leq \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) (\nu^{n+1} - \nu^{n}) (\mathrm{d}x) + L_{\nu} D_{h}(\nu^{n+1}, \nu^{n})$$

$$= -\frac{1}{\tau} \int_{\mathcal{X}} \left( \frac{\delta h}{\delta \nu} (\nu^{n+1}, x) - \frac{\delta h}{\delta \nu} (\nu^{n}, x) \right) (\nu^{n+1} - \nu^{n}) (\mathrm{d}x) + L_{\nu} D_{h}(\nu^{n+1}, \nu^{n})$$

$$= -\frac{1}{\tau} \left( D_{h}(\nu^{n+1}, \nu^{n}) + D_{h}(\nu^{n}, \nu^{n+1}) \right) + L_{\nu} D_{h}(\nu^{n+1}, \nu^{n}), \quad (28)$$

where the last equality follows from (21).

Similarly, using  $L_{\mu}$ -smoothness of  $\mu \mapsto F(\nu, \mu)$  relative to h together with (20), we can show that

$$F(\nu^{n+1},\mu^n) - F(\nu^{n+1},\mu^{n+1}) + \frac{1}{\tau} \left( D_h(\mu^n,\mu^{n+1}) + D_h(\mu^{n+1},\mu^n) \right) \le L_\mu D_h(\mu^{n+1},\mu^n). \tag{29}$$

Therefore, using (28) and (29) in (27), and recalling that  $L = \max\{L_{\nu}, L_{\mu}\}$  gives

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq \frac{1}{N\tau}\left(\sup_{\nu\in\mathcal{C}}D_{h}(\nu,\nu^{0}) + \sup_{\mu\in\mathcal{D}}D_{h}(\mu,\mu^{0})\right) + \frac{1}{2N\tau}\sum_{n=0}^{N-1}\left(D_{h}(\nu^{n},\nu^{n+1}) - D_{h}(\nu^{n+1},\nu^{n}) + D_{h}(\mu^{n},\mu^{n+1}) - D_{h}(\mu^{n+1},\mu^{n})\right) + \frac{L}{2N}\sum_{n=0}^{N-1}\left(D_{h}(\nu^{n+1},\nu^{n}) + D_{h}(\mu^{n+1},\mu^{n})\right) + \frac{1}{2N}\left(F\left(\nu^{N},\mu^{N}\right) - F(\nu^{0},\mu^{0})\right).$$
(30)

Since, by Lemma A.1,  $D_h(\nu^{n+1},\nu^n) \leq 4L_F\tau^2$  and  $D_h(\mu^{n+1},\mu^n) \leq 4L_F\tau^2$ , it suffices to show that  $D_h(\nu^n,\nu^{n+1}) - D_h(\nu^{n+1},\nu^n) + D_h(\mu^n,\mu^{n+1}) - D_h(\mu^{n+1},\mu^n)$  is of order  $\mathcal{O}(\tau^3)$ . Indeed, we could then choose  $\tau = \mathcal{O}\left(\frac{1}{N^{1/3}}\right)$ , and since by Assumption 3.5,  $\left|F\left(\nu^N,\mu^N\right)\right| \leq M$ , we would obtain that

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1},\frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq \mathcal{O}\left(\frac{1}{N^{2/3}}\right) + \mathcal{O}\left(\frac{1}{N}\right) = \mathcal{O}\left(\frac{1}{N^{2/3}}\right),$$

because  $\frac{1}{N} \leq \frac{1}{N^{2/3}}$ , for all  $N \geq 1$ .

 In order to show that  $D_h(\nu^n, \nu^{n+1}) - D_h(\nu^{n+1}, \nu^n) + D_h(\mu^n, \mu^{n+1}) - D_h(\mu^{n+1}, \mu^n)$  is  $\mathcal{O}(\tau^3)$ , we will leverage the connection between Bregman divergence and dual Bregman divergence given by Lemma H.3 together with Assumption 3.4, 3.5.

If we denote  $f^n := \frac{\delta h}{\delta \nu}(\nu^n, \cdot)$ , for any  $n \ge 0$ , then by Lemma H.3, we have that  $D_h(\nu^n, \nu^{n+1}) = D_{h^*}(f^{n+1}, f^n)$ . For any  $\varepsilon \in [0, 1]$  denote  $f^{\varepsilon, n+1} = \varepsilon f^{n+1} + (1 - \varepsilon) f^n$  and  $f^{\varepsilon, n} = \varepsilon f^n + (1 - \varepsilon) f^{n+1}$ . By Definition 3.3, we have that

$$\begin{split} D_{h^*}(f^{n+1},f^n) &= h^*(f^{n+1}) - h^*(f^n) - \int_{\mathcal{X}} \left(f^{n+1}(x) - f^n(x)\right) \frac{\delta h^*}{\delta f}(f^n)(\mathrm{d}x) \\ &= \int_0^1 \left\langle f^{n+1} - f^n, \frac{\delta h^*}{\delta f}(f^{\lambda,n+1}) \right\rangle \mathrm{d}\lambda - \left\langle f^{n+1} - f^n, \frac{\delta h^*}{\delta f}(f^n) \right\rangle \\ &= \int_0^1 \left\langle f^{n+1} - f^n, \frac{\delta h^*}{\delta f}(f^{\lambda,n+1}) - \frac{\delta h^*}{\delta f}(f^n) \right\rangle \mathrm{d}\lambda \\ &= \int_0^1 \int_0^1 \int_{\mathcal{X} \times \mathcal{X}} \lambda \left( f^{n+1}(x) - f^n(x) \right) \left( f^{n+1}(x') - f^n(x') \right) \frac{\delta^2 h^*}{\delta f^2}(f^{\eta\lambda,n+1})(f^{\eta\lambda,n+1})(\mathrm{d}x' \otimes \mathrm{d}x) \mathrm{d}\eta \mathrm{d}\lambda, \end{split}$$

where the second and last equalities follow from (43) and (46), respectively. Similarly, by Lemma H.3, we have that  $D_h(\nu^{n+1}, \nu^n) = D_{h^*}(f^n, f^{n+1})$ , and hence

$$D_{h^*}(f^n, f^{n+1}) = \int_0^1 \int_0^1 \int_{\mathcal{X} \times \mathcal{X}} \lambda \left( f^n(x) - f^{n+1}(x) \right) \left( f^n(x') - f^{n+1}(x') \right) \times \frac{\delta^2 h^*}{\delta f^2} (f^{\eta \lambda, n}) (f^{\eta \lambda, n}) (dx' \otimes dx) d\eta d\lambda.$$

Therefore, we obtain that

$$D_{h^*}(f^{n+1}, f^n) - D_{h^*}(f^n, f^{n+1}) = \int_0^1 \int_0^1 \int_{\mathcal{X} \times \mathcal{X}} \lambda(f^{n+1}(x) - f^n(x))(f^{n+1}(x') - f^n(x')) \times \left( \frac{\delta^2 h^*}{\delta f^2} (f^{\eta \lambda, n+1})(f^{\eta \lambda, n+1}) - \frac{\delta^2 h^*}{\delta f^2} (f^{\eta \lambda, n})(f^{\eta \lambda, n}) \right) (\mathrm{d}x' \otimes \mathrm{d}x) \mathrm{d}\eta \mathrm{d}\lambda.$$

Using Assumption 3.4, we further obtain

$$D_{h^*}(f^{n+1}, f^n) - D_{h^*}(f^n, f^{n+1}) \leq \int_0^1 \int_0^1 \left| \int_{\mathcal{X} \times \mathcal{X}} \lambda(f^{n+1}(x) - f^n(x))(f^{n+1}(x') - f^n(x')) \times \right| \\ \times \left( \frac{\delta^2 h^*}{\delta f^2} (f^{\eta \lambda, n+1})(f^{\eta \lambda, n+1}) - \frac{\delta^2 h^*}{\delta f^2} (f^{\eta \lambda, n})(f^{\eta \lambda, n}) \right) (\mathrm{d}x' \otimes \mathrm{d}x) \left| \mathrm{d}\eta \mathrm{d}\lambda \right| \\ \leq \|f^{n+1} - f^n\|_\infty^2 \int_0^1 \int_0^1 \int_{\mathcal{X} \times \mathcal{X}} \lambda \left| \frac{\delta^2 h^*}{\delta f^2} (f^{\eta \lambda, n+1})(f^{\eta \lambda, n+1}) - \frac{\delta^2 h^*}{\delta f^2} (f^{\eta \lambda, n})(f^{\eta \lambda, n}) \right| (\mathrm{d}x' \otimes \mathrm{d}x) \mathrm{d}\eta \mathrm{d}\lambda \\ = \|f^{n+1} - f^n\|_\infty^2 \int_0^1 \lambda \int_0^1 \left\| \frac{\delta^2 h^*}{\delta f^2} (f^{\eta \lambda, n+1})(f^{\eta \lambda, n+1}) - \frac{\delta^2 h^*}{\delta f^2} (f^{\eta \lambda, n})(f^{\eta \lambda, n}) \right\|_{\mathrm{TV}} \mathrm{d}\eta \mathrm{d}\lambda \\ \leq 2L_{h^*} \|f^{n+1} - f^n\|_\infty^3 \int_0^1 \lambda \int_0^1 |1 - 2\eta \lambda| \mathrm{d}\eta \mathrm{d}\lambda \leq L_{h^*} \|f^{n+1} - f^n\|_\infty^3,$$

where the third inequality follows since  $|1 - 2\eta\lambda| \le 1$ , for all  $\eta, \lambda \in [0, 1]$ . The first-order condition for the minimizing player in (20) can be rewritten as

$$f^{n+1}(x) - f^n(x) = -\tau \frac{\delta F}{\delta \nu}(\nu^n, \mu^n, x),$$
 (31)

for all  $x \in \mathcal{X}$  Lebesgue a.e. By Assumption 3.5, there exists  $C_{\nu} > 0$  such that  $\left\| \frac{\delta F}{\delta \nu}(\nu^n, \mu^n, \cdot) \right\|_{\infty} \le C_{\nu}$ , for any  $n \ge 0$ . Hence, we obtain that

$$D_{h^*}(f^{n+1}, f^n) - D_{h^*}(f^n, f^{n+1}) \le L_{h^*} \|f^{n+1} - f^n\|_{\infty}^3 = L_{h^*} \tau^3 \left\| \frac{\delta F}{\delta \nu}(\nu^n, \mu^n, \cdot) \right\|_{\infty}^3 \le L_{h^*} \tau^3 C_{\nu}^3,$$

Similarly, denoting  $g^n := \frac{\delta h}{\delta \mu}(\mu^n, \cdot)$ , for any  $n \geq 0$ , and repeating the steps above, we can prove that

$$D_{h^*}(g^{n+1}, g^n) - D_{h^*}(g^n, g^{n+1}) \le L_{h^*} \|g^{n+1} - g^n\|_{\infty}^3 = L_{h^*} \tau^3 \left\| \frac{\delta F}{\delta \mu}(\nu^{n+1}, \mu^n, \cdot) \right\|_{\infty}^3 \le L_{h^*} \tau^3 C_{\mu}^3,$$

where  $C_{\mu} > 0$  exists due Assumption 3.5.

Set  $\kappa := C_{\nu}^3 + C_{\mu}^3 > 0$ . Then,

$$D_h(\nu^n, \nu^{n+1}) - D_h(\nu^{n+1}, \nu^n) + D_h(\mu^n, \mu^{n+1}) - D_h(\mu^{n+1}, \mu^n) \le \kappa L_{h^*} \tau^3. \tag{32}$$

Hence, using Lemma A.1, (32) and Assumption 3.5, estimate (30) becomes

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq \frac{1}{N\tau}\left(\sup_{\nu\in\mathcal{C}}D_{h}(\nu,\nu^{0}) + \sup_{\mu\in\mathcal{D}}D_{h}(\mu,\mu^{0})\right) \\
+ \frac{1}{2N\tau}\sum_{n=0}^{N-1}\left(\left(D_{h}(\nu^{n},\nu^{n+1}) - D_{h}(\nu^{n+1},\nu^{n})\right) + \left(D_{h}(\mu^{n},\mu^{n+1}) - D_{h}(\mu^{n+1},\mu^{n})\right)\right) \\
+ \frac{L}{2N}\sum_{n=0}^{N-1}\left(D_{h}(\nu^{n+1},\nu^{n}) + D_{h}(\mu^{n+1},\mu^{n})\right) + \frac{1}{2N}\left(F\left(\nu^{N},\mu^{N}\right) - F(\nu^{0},\mu^{0})\right) \\
= \frac{1}{N\tau}\left(\sup_{\nu\in\mathcal{C}}D_{h}(\nu,\nu^{0}) + \sup_{\mu\in\mathcal{D}}D_{h}(\mu,\mu^{0})\right) + \left(\frac{1}{2}\kappa L_{h^{*}} + 4L_{F}L\right)\tau^{2} + \frac{M}{N}.$$

Minimizing the right-hand side over  $\tau$  amounts to taking

$$\tau = \left(\frac{\sup_{\nu \in \mathcal{C}} D_h(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}} D_h(\mu, \mu^0)}{N(\kappa L_{h^*} + 2L_F L)}\right)^{1/3},$$

and since  $\frac{1}{N} \leq \frac{1}{N^{2/3}}$ , for any  $N \geq 1$ , it follows that

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right) \leq \frac{1}{2N^{2/3}}\left(3\left(\sup_{\nu \in \mathcal{C}}D_{h}(\nu, \nu^{0}) + \sup_{\mu \in \mathcal{D}}D_{h}(\mu, \mu^{0})\right)^{2/3} \times \left(\kappa L_{h^{*}} + 2L_{F}L\right)\right)^{1/3} + 2M\right).$$

# B DETAILS ON NUMERICAL EXPERIMENTS

In this section, we present the additional details of the numerical experiments. We begin by summarizing the implementable versions of the simultaneous and sequential MDA algorithms introduced in Section 4. We now turn to Algorithms (3) and (4) in the setting where F corresponds to the GAN objective introduced in Example 1.2. Recall that F takes the form

$$F(\nu,\mu) = \int_{\mathcal{W}} \int_{\Theta} \int_{\mathcal{Y}} D_w(y) \left( T_{\theta} \# \xi - \hat{\xi} \right) (\mathrm{d}y) \nu(\mathrm{d}\theta) \mu(\mathrm{d}w)$$
$$= \int_{\mathcal{W}} \int_{\Theta} \int_{\mathcal{Y}} D_w(y) \left( T_{\theta} \# \xi \right) (\mathrm{d}y) \nu(\mathrm{d}\theta) \mu(\mathrm{d}w) - \int_{\mathcal{W}} \int_{\mathcal{Y}} D_w(y) \hat{\xi}(\mathrm{d}y) \mu(\mathrm{d}w)$$

# Algorithm 3: IMPLEMENTABLE SIMULTANEOUS MDA **Input:** objective function F, initial measures $(\nu^0, \mu^0)$ , stepsize $\tau, \gamma > 0$ , time horizons K, N and number of particles J

Generate i.i.d 
$$\left(X_{j}^{0}, Y_{j}^{0}\right)_{j=1}^{J} \sim (\nu^{0}, \mu^{0})$$
  
Set  $\left(X_{j,0}^{0}, Y_{j,0}^{0}\right)_{j=1}^{J} = \left(X_{j}^{0}, Y_{j}^{0}\right)_{j=1}^{J}$ 

for 
$$n = 0, 1, ..., N - 1$$
 do

for 
$$k = 0, 1, ..., K - 1$$
 do

Generate independent Gaussian random variables  $\mathcal{N}_{i,k}^n$ 

$$\begin{aligned} & \textbf{for } j = 1, 2, \dots, J \textbf{ do} \\ & & X_{j,k+1}^n = X_{j,k}^n - \gamma \nabla \frac{\delta F}{\delta \nu}(\nu^n, \mu^n, X_{j,k}^n) + \sqrt{\frac{2\gamma}{\tau}} \mathcal{N}_{j,k}^n \\ & & Y_{j,k+1}^n = Y_{j,k}^n + \gamma \nabla \frac{\delta F}{\delta \mu}(\nu^n, \mu^n, Y_{j,k}^n) + \sqrt{\frac{2\gamma}{\tau}} \mathcal{N}_{j,k}^n \end{aligned}$$

$$\begin{array}{l} \textbf{for } j = 1, 2, \dots, J \textbf{ do} \\ & \ \ \, \big\lfloor X_{j,0}^{n+1} = X_{j,K}^n, \quad Y_{j,0}^{n+1} = Y_{j,K}^n \\ & \ \ \, \nu^n = \frac{1}{J} \sum_{j=1}^J \delta_{X_{j,0}^n}, \quad \mu^n = \frac{1}{J} \sum_{j=1}^J \delta_{Y_{j,0}^n} \end{array}$$

**Output:** 
$$\left(\frac{1}{N} \sum_{n=0}^{N-1} \nu^n, \frac{1}{N} \sum_{n=0}^{N-1} \mu^n\right)$$

# Algorithm 4: IMPLEMENTABLE SEQUENTIAL MDA

**Input:** objective function F, initial measures  $(\nu^0, \mu^0)$ , stepsize  $\tau, \gamma > 0$ , time horizons K, N and number of particles J

Generate i.i.d 
$$\left(X_j^0,Y_j^0\right)_{j=1}^J \sim (\nu^0,\mu^0)$$

Set 
$$(X_{j,0}^0, Y_{j,0}^0)_{j=1}^J = (X_j^0, Y_j^0)_{j=1}^J$$

for k = 0, 1, ..., K - 1 do

for 
$$n = 0, 1, ..., N - 1$$
 do

Generate independent Gaussian random variables 
$$\mathcal{N}_{j,t}^n$$
 for  $j=1,2,\ldots,J$  do

$$\begin{bmatrix} X_{j,k+1}^n = X_{j,k}^n - \gamma \nabla \frac{\delta F}{\delta \nu}(\nu^n, \mu^n, X_{j,k}^n) + \sqrt{\frac{2\gamma}{\tau}} \mathcal{N}_{j,k}^n \end{bmatrix}$$

for 
$$j = 1, 2, ..., J$$
 do  
 $X_{i,0}^{n+1} = X_{i,K}^{n}$ 

$$u^{n+1} = \frac{1}{J} \sum_{i=1}^{J} \delta_{X_{i+1}^{n+1}}$$

for 
$$k = 0, 1, \dots, K - 1$$
 do

 $\begin{array}{ll} \textbf{for } k=0,1,\ldots,K-1 \textbf{ do} \\ | & \text{Generate independent Gaussian random variables } \mathcal{N}^n_{j,k} \end{array}$ 

$$\mu^n = \frac{1}{J} \sum_{j=1}^J \delta_{Y_{j,0}^n}$$

**Output:** 
$$\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n}\right)$$

By Definition F.1, we have

$$\frac{\delta F}{\delta \nu}(\nu, \mu, \theta) = \int_{\mathcal{W}} \int_{\mathcal{V}} D_w(y) \left( T_{\theta} \# \xi \right) (\mathrm{d}y) \mu(\mathrm{d}w),$$

$$\frac{\delta F}{\delta \mu}(\nu, \mu, w) = \int_{\Theta} \int_{\mathcal{V}} D_w(y) \left( T_{\theta} \# \xi \right) (\mathrm{d}y) \nu(\mathrm{d}\theta) - \int_{\mathcal{V}} D_w(y) \hat{\xi}(\mathrm{d}y).$$

The flat derivatives can be approximated using empirical averages. For a batch of real data  $\{\xi_1^{\text{real}}, \dots, \xi_M^{\text{real}}\} \sim \hat{\xi}$ , we have

$$\int_{\mathcal{Y}} D_w(y)\hat{\xi}(\mathrm{d}y) \approx \frac{1}{M} \sum_{i=1}^{M} D_w(\xi_i^{\text{real}}).$$

For the term in  $\frac{\delta F}{\delta \mu}(\nu, \mu, w)$  that involves integration with respect to both  $\nu$  and the generated data  $T_{\theta} \# \xi$ , we approximate via sampling as follows. We sample

$$\{\theta^{(1)}, \theta^{(2)}, ..., \theta^{(J)}\} \sim \nu, \quad \{Z_i^{(j)}\}_{i=1}^M \sim T_{\theta^{(j)}} \# \xi,$$

leading to the estimator

$$\int_{\Theta} \int_{\mathcal{Y}} D_w(y) \left( T_{\theta} \# \xi \right) (\mathrm{d}y) \nu(\mathrm{d}\theta) \approx \frac{1}{JM} \sum_{i=1}^{M} \sum_{j=1}^{J} D_w \left( X_i^{(j)} \right).$$

Analogously, for  $\frac{\delta F}{\delta \nu}(\nu, \mu, \theta)$  we sample

$$\{w^{(1)}, w^{(2)}, ..., w^{(J)}\} \sim \mu, \quad \{Z_i\}_{i=1}^M \sim T_\theta \# \xi,$$

and approximate

$$\int_{\mathcal{W}} \int_{\mathcal{Y}} D_w(y) \left( T_{\theta} \# \xi \right) (\mathrm{d}y) \mu(\mathrm{d}w) \approx \frac{1}{JM} \sum_{i=1}^{M} \sum_{j=1}^{J} D_{w^{(j)}} \left( Z_i \right).$$

To mitigate the computational cost of Algorithms (3) and (4), we follow the approach of (Hsieh et al., 2019) and employ Langevin dynamics with exponential damping (see also their Algorithm 3). Below, we present this algorithm in both the simultaneous and sequential variants used in our experiments.

# Algorithm 5: SIMULTANEOUS MDA-GAN

**Input:** Initial parameters  $w^0, \theta^0$ , step sizes  $\{\gamma^n\}_{n=0}^{N-1}, \{\tau^n\}_{n=0}^{N-1}$ , time horizon  $\{K^n\}_{n=0}^{N-1},$  averaging parameter  $\beta \in [0,1]$ , source probability measure  $\xi$ 

for 
$$n = 0, 1, \dots, N - 1$$
 do
$$\text{Set } \bar{w}^n, w_0^n = w^n \text{ and } \bar{\theta}^n, \theta_0^n = \theta^n;$$

$$\begin{aligned} & \text{for } k = 0, 1, \dots, K_t - 1 \text{ do} \\ & \mid A = \{Z_1, \dots, Z_M\} \sim T_{\theta_k^n} \# \xi; \end{aligned}$$

$$\theta_{k+1}^n = \theta_k^n - \frac{\gamma^n}{M} \nabla_{\theta} \sum_{Z_i \in A} D_{w^n}(Z_i) + \sqrt{\frac{2\gamma^n}{\tau^n}} \mathcal{N}_k^n;$$

$$B = \{\xi_1^{\text{real}}, \dots, \xi_M^{\text{real}}\} \sim \hat{\xi};$$
  

$$B' = \{Z'_1, \dots, Z'_M\} \sim T_{\theta^n} \# \xi;$$

$$w_{k+1}^n = w_k^n + \frac{\gamma^n}{M} \nabla_w \sum_{Z_i' \in B'} D_{w_k^n}(Z_i') - \frac{\gamma^n}{M} \nabla_w \sum_{\xi_i^{\text{real}} \in B} D_{w_k^n}(\xi_i^{\text{real}}) + \sqrt{\frac{2\gamma^n}{\tau^n}} \mathcal{N}_k^n;$$

$$\bar{w}^{n} = (1 - \beta)\bar{w}^{n} + \beta w_{k+1}^{n}, \quad \bar{\theta}^{n} = (1 - \beta)\bar{\theta}^{n} + \beta \theta_{k+1}^{n};$$

$$w^{n+1} = (1 - \beta)w^{t} + \beta \bar{w}^{n}, \quad \theta^{n+1} = (1 - \beta)\theta^{n} + \beta \bar{\theta}^{n};$$

Output:  $w^N, \theta^N$ 

In all experiments, we closely follow the specifications from (Hsieh et al., 2019). We adopt the gradient-penalized discriminator of (Gulrajani et al., 2017) as a soft-constraint alternative to the

#### Algorithm 6: SEQUENTIAL MDA-GAN

Output:  $w^N$ ,  $\theta^N$ 

```
 \begin{split} \overline{\textbf{Input:}} & \text{Initial parameters } w^0, \theta^0, \text{ step sizes } \{\gamma^n\}_{n=0}^{N-1}, \{\tau^n\}_{n=0}^{N-1}, \text{ time horizon } \{K^n\}_{n=0}^{N-1}, \\ & \text{averaging parameter } \beta \in [0,1], \text{ source probability measure } \xi \\ \hline \textbf{for } n = 0, 1, \dots, N-1 \textbf{ do} \\ & \text{Set } \bar{w}^n, w^n_0 = w^n \text{ and } \bar{\theta}^n, \theta^n_0 = \theta^n; \\ \textbf{ for } k = 0, 1, \dots, K_t-1 \textbf{ do} \\ & A = \{Z_1, \dots, Z_M\} \sim T_{\theta^n_k} \# \xi; \\ & \theta^n_{k+1} = \theta^n_k - \frac{\gamma^n}{M} \nabla_\theta \sum_{Z_i \in A} D_{w^n}(Z_i) + \sqrt{\frac{2\gamma^n}{\tau^n}} \mathcal{N}^n_k; \\ & \bar{\theta}^n = (1-\beta)\bar{\theta}^n + \beta\bar{\theta}^n; \\ \textbf{ for } k = 0, 1, \dots, K_t-1 \textbf{ do} \\ & B = \{\xi^{\text{real}}_1, \dots, \xi^{\text{real}}_M\} \sim \hat{\xi}; \\ & B' = \{Z'_1, \dots, Z'_M\} \sim T_{\theta^{n+1}} \# \xi; \\ & w^n_{k+1} = w^n_k + \frac{\gamma^n}{M} \nabla_w \sum_{Z'_i \in B'} D_{w^n_k}(Z'_i) - \frac{\gamma^n}{M} \nabla_w \sum_{\xi^{\text{real}}_i \in B} D_{w^n_k}(\xi^{\text{real}}_i) + \sqrt{\frac{2\gamma^n}{\tau^n}} \mathcal{N}^n_k; \\ & \bar{w}^n = (1-\beta)\bar{w}^n + \beta w^n_{k+1}; \\ & w^{n+1} = (1-\beta)w^t + \beta\bar{w}^n; \end{split}
```

original Wasserstein GAN formulation to increase stability. The gradient penalty parameter is set to  $\lambda=0.1$ . For our Simultaneous and Sequential MDA-GANs, we fix the damping factor to  $\beta=0.8$ . The scheduling of the parameters  $K^n, \gamma^n$ , and  $\tau^n$  is  $K^n=\lfloor (1+10^{-5})^n\rfloor, \gamma^n=\gamma(1-10^{-5})^n$ , with  $\gamma=0.01$ , and  $\tau^n=\tau(1-5\times10^{-5})^{-t}$ , with  $\tau=100$ . The number of samples per batch is M=1024. For both the 8-Gaussian mixture and Swiss Roll datasets, we use fully connected networks for the generator and discriminator, each consisting of two-hidden-layers with J=512 neurons on each layer. The generator and discriminator networks use ReLU activations, except for the output layer of the discriminator, which employs a tanh activation. All network parameters are initialized from a normal distribution  $\mathcal{N}(0,0.01)$ .

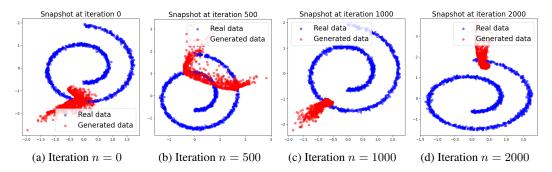


Figure 3: Simultaneous MDA-GAN (Algorithm 5) learning the Swiss Roll

# C PROOFS OF ADDITIONAL RESULTS

In this section, we present the proofs of the additional results of the paper. We start with the proofs of Lemma A.1 and Lemma C.1, which play a key role in proving the main results. Then we continue with the proofs of some auxiliary results.

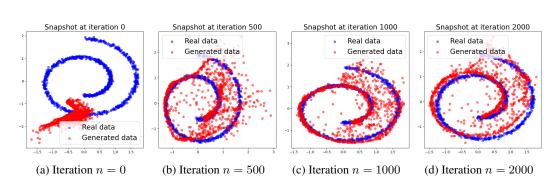


Figure 4: Sequential MDA-GAN (Algorithm 6) learning the Swiss Roll

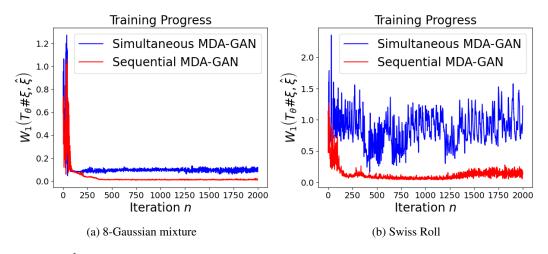


Figure 5:  $L^1$ -Wasserstein distance between generated and real data for the 8-Gaussian mixture and Swiss Roll

#### C.1 Proof of Lemma A.1

*Proof of Lemma A.1.* We will only prove the lemma for scheme (1) since the argument for (2) is almost identical. From  $L_{\nu}$ -relative smoothness and the definition of  $\nu^{n+1}$  in (1), for any  $\nu \in \mathcal{C}$ , it follows that

$$F(\nu^{n+1}, \mu^n) \leq F(\nu^n, \mu^n) + \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^n, \mu^n, x) (\nu^{n+1} - \nu^n) (\mathrm{d}x) + \left(\frac{1}{\tau} + L_{\nu} - \frac{1}{\tau}\right) D_h(\nu^{n+1}, \nu^n)$$

$$\leq F(\nu^n, \mu^n) + \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^n, \mu^n, x) (\nu - \nu^n) (\mathrm{d}x) + \frac{1}{\tau} D_h(\nu, \nu^n) + \left(L_{\nu} - \frac{1}{\tau}\right) D_h(\nu^{n+1}, \nu^n).$$

Setting  $\nu = \nu^n$ , we obtain that

$$F(\nu^{n+1}, \mu^n) \le F(\nu^n, \mu^n) + \left(L_{\nu} - \frac{1}{\tau}\right) D_h(\nu^{n+1}, \nu^n).$$

Recall  $L := \max\{L_{\nu}, L_{\mu}\} > 0$ . By assumption,  $\tau L \leq \frac{1}{2}$ , and so we get

$$\frac{1}{2\tau}D_h(\nu^{n+1},\nu^n) \le F(\nu^n,\mu^n) - F(\nu^{n+1},\mu^n) \le \sqrt{L_F}\sqrt{D_h(\nu^{n+1},\nu^n)},$$

where the last inequality follows from Assumption 2.1. Hence, since  $D_h(\nu^{n+1}, \nu^n) \ge 0$ , for all  $n \ge 0$ , we obtain that

$$D_h(\nu^{n+1}, \nu^n) \le 4L_F \tau^2.$$

From  $L_{\mu}$ -relative smoothness and the definition of  $\mu^{n+1}$  in (1), for any  $\mu \in \mathcal{D}$ , it follows that

$$F(\nu^{n}, \mu^{n+1}) \geq F(\nu^{n}, \mu^{n}) + \int_{\mathcal{X}} \frac{\delta F}{\delta \mu}(\nu^{n}, \mu^{n}, y)(\mu^{n+1} - \mu^{n})(\mathrm{d}y) - \left(\frac{1}{\tau} + L_{\mu} - \frac{1}{\tau}\right) D_{h}(\mu^{n+1}, \mu^{n})$$

$$\geq F(\nu^{n}, \mu^{n}) + \int_{\mathcal{X}} \frac{\delta F}{\delta \mu}(\nu^{n}, \mu^{n}, y)(\mu - \mu^{n})(\mathrm{d}y) - \frac{1}{\tau} D_{h}(\mu, \mu^{n}) - \left(L_{\mu} - \frac{1}{\tau}\right) D_{h}(\mu^{n+1}, \mu^{n}).$$

Setting  $\mu = \mu^n$ , we obtain that

$$F(\nu^n, \mu^{n+1}) \ge F(\nu^n, \mu^n) - \left(L_\mu - \frac{1}{\tau}\right) D_h(\mu^{n+1}, \mu^n).$$

Using again the assumption  $\tau L \leq \frac{1}{2}$ , we get

$$\frac{1}{2\tau}D_h(\mu^{n+1},\mu^n) \le F(\nu^n,\mu^{n+1}) - F(\nu^n,\mu^n) \le \sqrt{L_F}\sqrt{D_h(\mu^{n+1},\mu^n)}$$

where the last inequality follows from Assumption 2.1. Hence, since  $D_h(\mu^{n+1}, \mu^n) \ge 0$ , for all  $n \ge 0$ , we obtain that

$$D_h(\mu^{n+1}, \mu^n) \le 4L_F \tau^2.$$

#### C.2 Proof of Lemma C.1

**Lemma C.1** (Three-point inequality). Let Assumption 1.1 hold. Let  $G : \mathcal{E} \to \mathbb{R}$  be convex and admit flat derivative on  $\mathcal{E}$ . For all  $\mu \in \mathcal{E}$ , suppose that there exists  $\bar{\nu} \in \mathcal{E}$  such that

$$\bar{\nu} \in \underset{\nu \in \mathcal{E}}{\operatorname{arg\,min}} \{ G(\nu) + D_h(\nu, \mu) \}.$$

*Then, for any*  $\nu \in \mathcal{E}$ *, we have* 

$$G(\nu) + D_h(\nu, \mu) \ge G(\bar{\nu}) + D_h(\bar{\nu}, \mu) + D_h(\nu, \bar{\nu}).$$

*Proof.* From Definition 1.2, we have

$$D_h(\nu,\mu) = h(\nu) - h(\mu) - \int_{\mathcal{V}} \frac{\delta h}{\delta \mu}(\mu,y)(\nu - \mu)(\mathrm{d}y),$$

and hence, for any  $\mu \in \mathcal{K}$ , and all  $y \in \mathcal{X}$  Lebesgue a.e., we have

$$\left(\frac{\delta D_h}{\delta \nu}(\nu, \mu, y)\right)\bigg|_{\nu = \bar{\nu}} = \frac{\delta h}{\delta \nu}(\bar{\nu}, y) - \frac{\delta h}{\delta \mu}(\mu, y).$$

Therefore, for any  $\mu \in \mathcal{K}$ , we have that

$$D_{D_{h}(\cdot,\mu)}(\nu,\bar{\nu}) = D_{h}(\nu,\mu) - D_{h}(\bar{\nu},\mu) - \int_{\mathcal{X}} \left( \frac{\delta D_{h}}{\delta \nu}(\nu,\mu,y) \right) \Big|_{\nu=\bar{\nu}} (\nu-\bar{\nu})(\mathrm{d}y)$$

$$= D_{h}(\nu,\mu) - D_{h}(\bar{\nu},\mu) - \int_{\mathcal{X}} \frac{\delta h}{\delta \nu}(\bar{\nu},y)(\nu-\bar{\nu})(\mathrm{d}y) + \int_{\mathcal{X}} \frac{\delta h}{\delta \mu}(\mu,y)(\nu-\bar{\nu})(\mathrm{d}y)$$

$$= h(\nu) - h(\bar{\nu}) - \int_{\mathcal{X}} \frac{\delta h}{\delta \mu}(\mu,y)(\nu-\mu)(\mathrm{d}y) + \int_{\mathcal{X}} \frac{\delta h}{\delta \mu}(\mu,y)(\bar{\nu}-\mu)(\mathrm{d}y)$$

$$- \int_{\mathcal{X}} \frac{\delta h}{\delta \nu}(\bar{\nu},y)(\nu-\bar{\nu})(\mathrm{d}y) + \int_{\mathcal{X}} \frac{\delta h}{\delta \mu}(\mu,y)(\nu-\bar{\nu})(\mathrm{d}y)$$

$$= h(\nu) - h(\bar{\nu}) - \int_{\mathcal{X}} \frac{\delta h}{\delta \nu}(\bar{\nu},y)(\nu-\bar{\nu})(\mathrm{d}y)$$

$$= D_{h}(\nu,\bar{\nu}).$$

Given  $\mu \in \mathcal{K}$ , if we denote  $g(\nu) \coloneqq G(\nu) + D_h(\nu, \mu)$ , then by linearity of flat derivative, we further obtain that

$$D_{g}(\nu,\bar{\nu}) = D_{G(\cdot)+D_{h}(\cdot,\mu)}(\nu,\bar{\nu}) = D_{G}(\nu,\bar{\nu}) + D_{D_{h}(\cdot,\mu)}(\nu,\bar{\nu}) = D_{G}(\nu,\bar{\nu}) + D_{h}(\nu,\bar{\nu}) \ge D_{h}(\nu,\bar{\nu}),$$

since  $D_G(\nu, \bar{\nu}) \geq 0$  by convexity of G. By optimality of  $\bar{\nu}$ , the first-order condition  $\frac{\delta g}{\delta \nu}(\bar{\nu}, y) = \text{constant holds for all } y \in \mathcal{X}$  Lebesgue a.e., and hence

$$g(\nu) - g(\bar{\nu}) - D_g(\nu, \bar{\nu}) = 0.$$

Therefore, we obtain that

$$g(\nu) = g(\bar{\nu}) + D_g(\nu, \bar{\nu}) \ge g(\bar{\nu}) + D_h(\nu, \bar{\nu}),$$

which is the desired inequality.

#### C.3 PROOFS OF AUXILIARY RESULTS

In this subsection, we start by establishing two results: one concerning the verification of Assumption 2.1 and the other on the uniform boundedness of the second-order flat derivatives of F.

**Lemma C.2** (Verification of Assumption 2.1 for h relative entropy). Suppose that there exists  $C_{F,\nu}>0$  and  $C_{F,\mu}>0$  such that, for all  $(\nu,\mu)\in\mathcal{C}\times\mathcal{D}$ , and all  $(x,y)\in\mathcal{X}\times\mathcal{X}$ , it holds that

$$\left| \frac{\delta F}{\delta \nu}(\nu, \mu, x) \right| \le C_{F, \nu}, \left| \frac{\delta F}{\delta \mu}(\nu, \mu, y) \right| \le C_{F, \mu}.$$

Take h to be the relative entropy, i.e.,  $h(\nu) \coloneqq \int_{\mathcal{X}} \log \frac{\nu(x)}{\pi(x)} \nu(\mathrm{d}x)$ , where  $\nu, \pi \in \mathcal{P}(\mathcal{X})$  are absolutely continuous with respect to Lebesgue measure on  $\mathcal{X}$  and  $\pi$  is fixed reference probability measures on  $\mathcal{P}(\mathcal{X})$ . Then Assumption 2.1 is satisfied.

*Proof.* Since h is the relative entropy, it follows from Example 1.3 that Bregman divergence is in fact the Kullback-Leibler divergence. Then, from Definition F.1, we have

$$|F(\nu', \mu') - F(\nu, \mu)| = |F(\nu', \mu') - F(\nu, \mu') + F(\nu, \mu') - F(\nu, \mu)|$$

$$\leq |F(\nu', \mu') - F(\nu, \mu')| + |F(\nu, \mu') - F(\nu, \mu)|$$

$$= \left| \int_0^1 \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} \left( \nu + \varepsilon(\nu' - \nu), \mu', x \right) (\nu' - \nu) (\mathrm{d}x) \mathrm{d}\varepsilon \right|$$

$$+ \left| \int_0^1 \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} \left( \nu, \mu + \varepsilon(\mu' - \mu), y \right) (\mu' - \mu) (\mathrm{d}y) \mathrm{d}\varepsilon \right|$$

$$\leq C_{F,\nu} \operatorname{TV} \left( \nu', \nu \right) + C_{F,\mu} \operatorname{TV} \left( \mu', \mu \right),$$

where the last inequality follows since F is assumed to have bounded first-order flat derivatives by  $C_{F,\nu}, C_{F,\mu} > 0$ , respectively. The conclusion follows by squaring both sides and applying Pinsker's inequality, that is,  $\mathrm{TV}^2(\nu',\nu) \leq \frac{1}{2} \mathrm{KL}(\nu',\nu)$ .

We show that, under Assumption 1.5 and 1.6, the second-order flat derivatives  $\frac{\delta^2 F}{\delta \nu^2}$ ,  $-\frac{\delta^2 F}{\delta \mu^2}$  are nonnegative and bounded above by  $\frac{\delta^2 h}{\delta \nu^2}$ ,  $\frac{\delta^2 h}{\delta \mu^2}$  multiplied by the respective smoothness constants.

**Lemma C.3** (Uniform boundedness of second order flat derivatives of F). Let Assumption 1.5, 1.1 and 1.6 hold. Suppose that  $\nu \mapsto F(\nu, \mu)$ ,  $\mu \mapsto F(\nu, \mu)$ , and h admit second-order flat derivative (cf. (40)) on  $\mathcal{C}, \mathcal{D}$  and  $\mathcal{E}$ , respectively. Then, we have

$$0 \leq \int_{0}^{1} \int_{\mathcal{X}} \int_{0}^{\varepsilon} \int_{\mathcal{X}} \frac{\delta^{2} F}{\delta \nu^{2}} \left(\nu + \eta(\nu' - \nu), \mu, x, x'\right) (\nu' - \nu) (\mathrm{d}x') \mathrm{d}\eta (\nu' - \nu) (\mathrm{d}x) \mathrm{d}\varepsilon$$
$$\leq L_{\nu} \int_{0}^{1} \int_{\mathcal{X}} \int_{0}^{\varepsilon} \int_{\mathcal{X}} \frac{\delta^{2} h}{\delta \nu^{2}} \left(\nu + \eta(\nu' - \nu), x, x'\right) (\nu' - \nu) (\mathrm{d}x') \mathrm{d}\eta (\nu' - \nu) (\mathrm{d}x) \mathrm{d}\varepsilon,$$

$$0 \leq -\int_{0}^{1} \int_{\mathcal{X}} \int_{0}^{\varepsilon} \int_{\mathcal{X}} \frac{\delta^{2} F}{\delta \mu^{2}} \left(\nu, \mu + \eta(\mu' - \mu), y, y'\right) (\mu' - \mu) (\mathrm{d}y') \mathrm{d}\eta (\mu' - \mu) (\mathrm{d}y) \mathrm{d}\varepsilon$$
$$\leq L_{\mu} \int_{0}^{1} \int_{\mathcal{X}} \int_{0}^{\varepsilon} \int_{\mathcal{X}} \frac{\delta^{2} h}{\delta \mu^{2}} \left(\mu + \eta(\mu' - \mu), y, y'\right) (\mu' - \mu) (\mathrm{d}y') \mathrm{d}\eta (\mu' - \mu) (\mathrm{d}y) \mathrm{d}\varepsilon.$$

*Proof.* We observe that combining relative smoothness and convexity for  $\nu \mapsto F(\nu, \mu)$  gives that for some  $L_{\nu} > 0$ , any  $\nu, \nu' \in \mathcal{C}$  and any  $\mu, \mu' \in \mathcal{D}$ , we have

$$0 \le F(\nu', \mu) - F(\nu, \mu) - \int_{\mathcal{X}} \frac{\delta F}{\delta \nu}(\nu, \mu, x)(\nu' - \nu)(\mathrm{d}x) \le L_{\nu} D_h(\nu', \nu). \tag{33}$$

Since  $\nu \mapsto F(\nu, \mu)$ ,  $\mu \mapsto F(\nu, \mu)$ , and h admit second-order flat derivative (cf. (40)) on  $\mathcal{C}, \mathcal{D}$  and  $\mathcal{E}$ , respectively, from (33), we obtain

$$0 \leq \int_{0}^{1} \int_{\mathcal{X}} \int_{0}^{\varepsilon} \int_{\mathcal{X}} \frac{\delta^{2} F}{\delta \nu^{2}} \left(\nu + \eta(\nu' - \nu), \mu, x, x'\right) (\nu' - \nu) (\mathrm{d}x') \mathrm{d}\eta (\nu' - \nu) (\mathrm{d}x) \mathrm{d}\varepsilon$$
$$\leq L_{\nu} \int_{0}^{1} \int_{\mathcal{X}} \int_{0}^{\varepsilon} \int_{\mathcal{X}} \frac{\delta^{2} h}{\delta \nu^{2}} \left(\nu + \eta(\nu' - \nu), x, x'\right) (\nu' - \nu) (\mathrm{d}x') \mathrm{d}\eta (\nu' - \nu) (\mathrm{d}x) \mathrm{d}\varepsilon.$$

The analogous inequalities are similarly obtained for relative smoothness and relative concavity.  $\Box$ 

When F is strongly-convex-strongly-concave relative to h and Assumption 1.1 holds, it can be shown that  $(\nu^*, \mu^*)$  is the unique MNE of (1) (see the proof of (Lascu et al., 2025, Lemma A.5)). Moreover, based on relative convexity-concavity of F, we prove in Lemma C.5 that the NI error satisfies a type of "quadratic growth" inequality relative to h.

**Assumption C.4** (Relative convexity-concavity). Assume that, given  $\ell_{\nu}, \ell_{\mu} > 0$ , the function F is  $\ell_{\nu}$ -strongly convex in  $\nu$  and  $\ell_{\mu}$ -strongly concave in  $\mu$  relative to h, i.e., for any  $\nu, \nu' \in \mathcal{C}$  and any  $\mu, \mu' \in \mathcal{D}$ , we have

$$D_{F(\cdot,\mu)}(\nu',\nu) = F(\nu',\mu) - F(\nu,\mu) - \int_{\mathcal{X}} \frac{\delta F}{\delta \nu}(\nu,\mu,x)(\nu'-\nu)(\mathrm{d}x) \ge \ell_{\nu} D_{h}(\nu',\nu), \tag{34}$$

$$D_{F(\nu,\cdot)}(\mu',\mu) = F(\nu,\mu') - F(\nu,\mu) - \int_{\mathcal{V}} \frac{\delta F}{\delta \mu}(\nu,\mu,y)(\mu'-\mu)(\mathrm{d}y) \le -\ell_{\mu} D_{h}(\mu',\mu). \tag{35}$$

**Lemma C.5** ("Quadratic growth" of NI error relative to h). Suppose that Assumption 1.1 and C.4 hold. Then, for all  $(\nu, \mu) \in \mathcal{C} \times \mathcal{D}$ , it holds that

$$NI(\nu, \mu) \ge \ell (D_h(\nu, \nu^*) + D_h(\mu, \mu^*)),$$

where  $\ell := \min\{\ell_{\nu}, \ell_{\mu}\}.$ 

**Remark C.6.** We refer to Lemma C.5 as "quadratic growth" of NI error relative to h due to the similar notion of quadratic growth of a convex function relative to the squared Euclidean norm on  $\mathbb{R}^d$  (see e.g. (Anitescu, 2000)).

*Proof.* Let  $(\nu, \mu) \in \mathcal{C} \times \mathcal{D}$ . Since F is  $\ell_{\nu}$ -strongly convex in  $\nu$  and  $\ell_{\mu}$ -strongly concave in  $\mu$ , it follows that

$$F(\nu, \mu^*) - F(\nu^*, \mu^*) \ge \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^*, \mu^*, x) (\nu - \nu^*) (\mathrm{d}x) + \ell_{\nu} D_h(\nu, \nu^*),$$

$$F(\nu^*, \mu) - F(\nu^*, \mu^*) \le \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^*, \mu^*, y) (\mu - \mu^*) (\mathrm{d}y) - \ell_{\mu} D_h(\mu, \mu^*).$$

Since  $(\nu^*, \mu^*)$  is the MNE of F, we have

$$\frac{\delta F}{\delta \nu}(\nu^*,\mu^*,x) = \text{constant}, \quad \frac{\delta F}{\delta \mu}(\nu^*,\mu^*,y) = \text{constant},$$

for all  $(x,y) \in \mathcal{X} \times \mathcal{X}$  Lebesgue a.e. Hence, adding the inequalities above and using the definition of NI error, we get

$$NI(\nu,\mu) \ge \ell \left( D_h(\nu,\nu^*) + D_h(\mu,\mu^*) \right).$$

By Lemma C.5, the time-averaged iterates  $\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^n,\frac{1}{N}\sum_{n=0}^{N-1}\mu^n\right)$  converge in Bregman divergence to the unique MNE  $(\nu^*,\mu^*)$  of (1) with the rates proved in Theorem 2.4 and Theorem 3.6, respectively.

We now check that the condition  $\sup_{\nu \in \mathcal{C}} D_h(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}} D_h(\mu, \mu^0) < \infty$  required in Theorems 2.4 and 3.6 is satisfied in the specific cases of Examples 1.3 and 1.4.

**Lemma C.7.** Let h denote the relative entropy from Example 1.3, and set  $\mathcal{E}_{\beta} = \mathcal{C} \cup \mathcal{D}$ . Suppose  $\nu_0, \mu_0 \in \mathcal{E}_{\beta}$ , and assume there exists  $C_1, C_2 > 0$  such that, for all  $\nu, \mu \in \mathcal{P}(\mathcal{X})$ ,

$$\left\| \frac{\delta F}{\delta \nu}(\nu, \mu, \cdot) \right\|_{L^{\infty}(\mathcal{X})} \le C_1, \quad \left\| \frac{\delta F}{\delta \mu}(\nu, \mu, \cdot) \right\|_{L^{\infty}(\mathcal{X})} \le C_2.$$

Then the iterates produced by Algorithms (1) and (2) remain in  $\mathcal{E}_{\beta}$ , i.e.,

$$(\nu^n,\mu^n)_{n\geq 0}\subset \mathcal{E}_{\beta}.$$

Furthermore, they satisfy the uniform bound

$$\sup_{\nu \in \mathcal{C}} \mathrm{KL}(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}} \mathrm{KL}(\mu, \mu^0) \le 6\beta + 2\tau C_1 + 2\tau C_2.$$

*Proof.* We provide the proof only for Algorithm (1), as the argument for the other algorithm is essentially the same. Using the flat derivative formula (4), the first-order optimality condition (Hu et al., 2021, Proposition 2.5) applied to  $(\nu^{n+1}, \mu^{n+1})$  in Algorithm (1) gives

$$\begin{cases} \log \frac{\nu^{1}(x)}{\pi(x)} - \log \frac{\nu^{0}(x)}{\pi(x)} = -\tau \frac{\delta F}{\delta \nu}(\nu^{0}, \mu^{0}, x) - \log \int_{\mathcal{X}} e^{-\tau \frac{\delta F}{\delta \nu}(\nu^{0}, \mu^{0}, x)} \frac{\nu^{0}(x)}{\pi(x)} \pi(x) dx, \\ \log \frac{\mu^{1}(y)}{\pi(y)} - \log \frac{\mu^{0}(y)}{\pi(y)} = \tau \frac{\delta F}{\delta \mu}(\nu^{0}, \mu^{0}, y) - \log \int_{\mathcal{X}} e^{\tau \frac{\delta F}{\delta \mu}(\nu^{0}, \mu^{0}, y)} \frac{\mu^{0}(y)}{\pi(y)} \pi(y) dy, \end{cases}$$

for all  $x, y \in \mathcal{X}$  a.e. Taking the sup-norm on both sides over x, y and using the assumptions gives

$$\left\| \log \frac{\nu^1(\cdot)}{\pi(\cdot)} \right\|_{L^{\infty}(\mathcal{X})} \le 2\beta + 2\tau C_1,$$

$$\left\| \log \frac{\mu^1(\cdot)}{\pi(\cdot)} \right\|_{L^{\infty}(\mathcal{X})} \le 2\beta + 2\tau C_2,$$

and inductively,  $(\nu^n, \mu^n)_{n>0} \subset \mathcal{E}_{\beta}$ . Therefore, for any  $(\nu, \mu) \in \mathcal{E}_{\beta}$ ,

$$KL(\nu, \nu^{0}) + KL(\mu, \mu^{0}) = \int_{\mathcal{X}} \left( \log \frac{\nu(x)}{\pi(x)} - \log \frac{\nu^{0}(x)}{\pi(x)} \right) \nu(x) dx$$
$$+ \int_{\mathcal{X}} \left( \log \frac{\mu(y)}{\pi(y)} - \log \frac{\mu^{0}(y)}{\pi(y)} \right) \mu(y) dy$$
$$\leq 6\beta + 2\tau C_{1} + 2\tau C_{2},$$

hence the conclusion.

 **Lemma C.8.** Let h denote the  $\chi^2$ -divergence from Example 1.4, and set  $\mathcal{F}_{\eta} = \mathcal{C} \cup \mathcal{D}$ . Suppose  $\nu_0, \mu_0 \in \mathcal{F}_{\eta}$ , and assume there exists  $C_1, C_2 > 0$  such that, for all  $\nu, \mu \in \mathcal{P}(\mathcal{X})$ ,

$$\left\| \frac{\delta F}{\delta \nu}(\nu, \mu, \cdot) \right\|_{L^2_{\pi}(\mathcal{X})} \le C_1, \quad \left\| \frac{\delta F}{\delta \mu}(\nu, \mu, \cdot) \right\|_{L^2_{\pi}(\mathcal{X})} \le C_2.$$

Then the iterates produced by Algorithms (1) and (2) remain in  $\mathcal{F}_{\eta}$ , i.e.,

$$(\nu^n,\mu^n)_{n\geq 0}\subset \mathcal{F}_{\eta}.$$

Furthermore, they satisfy the uniform bound

$$\frac{1}{2} \sup_{\nu \in \mathcal{C}} \left\| \frac{\nu(\cdot)}{\pi(\cdot)} - \frac{\nu^0(\cdot)}{\pi(\cdot)} \right\|_{L^2_{\pi}(\mathcal{X})}^2 + \frac{1}{2} \sup_{\mu \in \mathcal{D}} \left\| \frac{\mu(\cdot)}{\pi(\cdot)} - \frac{\mu^0(\cdot)}{\pi(\cdot)} \right\|_{L^2_{\pi}(\mathcal{X})}^2 \le 4\eta + \tau C_1 + \tau C_2.$$

*Proof.* We provide the proof only for Algorithm (1), as the argument for the other algorithm is essentially the same. The first-order condition (see e.g., (Bonnans & Shapiro, 2000, Section 5.1.1)) shows that for a.e.  $x, y \in \mathcal{X}$ ,

$$\begin{split} &\left\langle \frac{\delta F}{\delta \nu}(\nu^0, \mu^0, \cdot) + \frac{1}{\tau} \left( \frac{\mathrm{d} \nu^1}{\mathrm{d} \pi} - \frac{\mathrm{d} \nu^0}{\mathrm{d} \pi} \right), \phi - \frac{\mathrm{d} \nu^1}{\mathrm{d} \pi} \right\rangle_{L^2_\pi} \geq 0, \quad \forall \phi \in \mathfrak{C}\,, \\ &\left\langle \frac{\delta F}{\delta \mu}(\nu^0, \mu^0, \cdot) - \frac{1}{\tau} \left( \frac{\mathrm{d} \mu^1}{\mathrm{d} \pi} - \frac{\mathrm{d} \mu^0}{\mathrm{d} \pi} \right), \phi - \frac{\mathrm{d} \mu^1}{\mathrm{d} \pi} \right\rangle_{L^2} \geq 0, \quad \forall \phi \in \mathfrak{C}\,, \end{split}$$

where  $\langle \cdot, \cdot \rangle_{L^2_{\pi}}$  is the inner product on  $L^2_{\pi}(\mathcal{X})$ , and  $\mathfrak{C}$  is the nonempty closed convex set defined by

$$\mathfrak{C} = \left\{\phi \in L^2_\pi(\mathcal{X}) \middle| \phi \geq 0 \text{ $\pi$-a.e. on $\mathcal{X}$ and } \int \phi(x) \pi(\mathrm{d}x) = 1\right\}.$$

Define the projection map  $\Pi_{\mathfrak{C}}: L^2_{\pi}(\mathcal{X}) \mapsto \mathfrak{C}$  such that  $\Pi_{\mathfrak{C}}(\varphi) = \arg\min_{\phi \in \mathfrak{C}} \|\phi - \varphi\|_{L^2_{\pi}(\mathcal{X})}$  for all  $\varphi \in L^2_{\pi}(\mathcal{X})$ , which satisfies

$$\langle \Pi(\varphi) - \varphi, \phi - \Pi(\varphi) \rangle_{L^2} \ge 0, \quad \forall \phi \in \mathfrak{C}.$$

Then

$$\begin{split} \frac{\mathrm{d}\nu^1}{\mathrm{d}\pi} &= \Pi_{\mathfrak{C}} \left( \frac{\mathrm{d}\nu^0}{\mathrm{d}\varrho} - \tau \frac{\delta F}{\delta \nu} (\nu^0, \mu^0, \cdot) \right), \\ \frac{\mathrm{d}\mu^1}{\mathrm{d}\pi} &= \Pi_{\mathfrak{C}} \left( \frac{\mathrm{d}\mu^0}{\mathrm{d}\varrho} + \tau \frac{\delta F}{\delta \mu} (\nu^0, \mu^0, \cdot) \right). \end{split}$$

Note  $\|\Pi_{\mathfrak{C}}(\varphi_1) - \Pi_{\mathfrak{C}}(\varphi_2)\|_{L^2_{\pi}(\mathcal{X})} \leq \|\varphi_1 - \varphi_2\|_{L^2_{\pi}(\mathcal{X})}$  for all  $\varphi_1, \varphi_2 \in L^2_{\pi}(\mathcal{X})$  (see e.g., (Ciarlet, 2013, Theorem 4.3-1)). Moreover, since  $\frac{\mathrm{d}\nu^0}{\mathrm{d}\pi} = \Pi_{\mathfrak{C}}\left(\frac{\mathrm{d}\nu^0}{\mathrm{d}\pi}\right)$ ,  $\frac{\mathrm{d}\mu^0}{\mathrm{d}\pi} = \Pi_{\mathfrak{C}}\left(\frac{\mathrm{d}\mu^0}{\mathrm{d}\pi}\right)$ , for a.e.  $x, y \in \mathcal{X}$ ,

$$\left\| \frac{\mathrm{d}\nu^{1}}{\mathrm{d}\pi} \right\|_{L_{\pi}^{2}(A)} \leq \left\| \frac{\mathrm{d}\nu^{0}}{\mathrm{d}\pi} \right\|_{L_{\pi}^{2}(A)} + \tau \left\| \frac{\delta F}{\delta \nu} (\nu^{0}, \mu^{0}, \cdot) \right\|_{L_{\pi}^{2}(A)} \leq \eta + \tau C_{1},$$

$$\left\| \frac{\mathrm{d}\mu^{1}}{\mathrm{d}\pi} \right\|_{L^{2}(A)} \leq \left\| \frac{\mathrm{d}\mu^{0}}{\mathrm{d}\pi} \right\|_{L^{2}(A)} + \tau \left\| \frac{\delta F}{\delta \mu} (\nu^{0}, \mu^{0}, \cdot) \right\|_{L^{2}(A)} \leq \eta + \tau C_{2},$$

and inductively,  $(\nu^n, \mu^n)_{n\geq 0} \subset \mathcal{F}_{\eta}$ . Therefore, for any  $(\nu, \mu) \in \mathcal{F}_{\eta}$ ,

$$\frac{1}{2} \left\| \frac{\nu(\cdot)}{\pi(\cdot)} - \frac{\nu^0(\cdot)}{\pi(\cdot)} \right\|_{L^2(\mathcal{X})}^2 + \frac{1}{2} \left\| \frac{\mu(\cdot)}{\pi(\cdot)} - \frac{\mu^0(\cdot)}{\pi(\cdot)} \right\|_{L^2(\mathcal{X})}^2 \le 4\eta + \tau C_1 + \tau C_2,$$

hence the conclusion.  $\Box$ 

# D VERIFICATION OF ASSUMPTION 1.5, 1.6, 2.1 AND 3.5 FOR EXAMPLE 1.2

In this section we verify that Assumption 1.5, 1.6, 2.1, and 3.5 are satisfied by the objective function F in Example 1.2.

**Proposition D.1** (Verification of assumptions for Example 1.2). Let  $\mathcal{Y}, \mathcal{Z} \subset \mathbb{R}^d$ , with  $\hat{\xi} \in \mathcal{P}(\mathcal{Y})$  and  $\xi \in \mathcal{P}(\mathcal{Z})$ . Suppose  $T_{\theta} : \mathcal{Z} \to \mathcal{Y}$  is measurable with  $\theta \in \Theta \subset \mathbb{R}^d$ , and  $D_w : \mathcal{Y} \to \mathbb{R}$  is uniformly bounded and measurable with  $w \in \mathcal{W} \subset \mathbb{R}^d$ . Then Assumptions 1.5, 1.6, 2.1 and 3.5 are satisfied by the objective

$$F(\nu,\mu) := \int_{\mathcal{W}} \int_{\Theta} f(\theta,w) \nu(\mathrm{d}\theta) \mu(\mathrm{d}w)$$

from Example 1.2.

*Proof.* By Definition F.1,

$$\frac{\delta F}{\delta \nu}(\nu,\mu,\theta) = \int_{\mathcal{W}} f(\theta,w) \mu(\mathrm{d}w),$$

and

$$\frac{\delta F}{\delta \mu}(\nu, \mu, w) = \int_{\Theta} f(\theta, w) \nu(\mathrm{d}\theta).$$

Therefore, Assumption 1.5 holds with equality, and Assumption 1.6 holds with equality with  $L_{\nu}=L_{\mu}=0$ . Since  $D_{w}$  is uniformly bounded by some  $M_{D}>0$ , we have

$$|f(\theta, w)| = \left| \int_{\mathcal{Y}} D_w(y) \left( T_{\theta} \# \xi - \hat{\xi} \right) (\mathrm{d}y) \right|$$

$$\leq \int_{\mathcal{Y}} |D_w(y)| \left( T_{\theta} \# \xi \right) (\mathrm{d}y) + \int_{\mathcal{Y}} |D_w(y)| \hat{\xi}(\mathrm{d}y) \leq 2M_D,$$

where the last inequality holds because  $\hat{\xi} \in \mathcal{P}(\mathcal{Y})$  and, since  $\xi \in \mathcal{P}(\mathcal{Z})$ , we have  $T_{\theta} \# \xi \in \mathcal{P}(\mathcal{Y})$ . Therefore,

$$\begin{split} |F(\nu',\mu') - F(\nu,\mu)| &= |F(\nu',\mu') - F(\nu',\mu) + F(\nu',\mu) - F(\nu,\mu)| \\ &\leq \int_w \int_\theta |f(\theta,w)| \nu'(\mathrm{d}\theta) |\mu' - \mu| (\mathrm{d}w) \\ &+ \int_w \int_\theta |f(\theta,w)| |\nu' - \nu| (\mathrm{d}\theta) \mu(\mathrm{d}w) \\ &\leq 4 M_D \left( \mathrm{TV}(\nu',\nu) + \mathrm{TV}(\mu',\mu) \right), \end{split}$$

since  $\mathrm{TV}(m',m)=\frac{1}{2}\int |m'-m|\mathrm{d}x$  for all  $m\in\mathcal{P}(\mathcal{X})$  by (Tsybakov, 2008, Lemma 2.1). Squaring both sides and applying Remark 2.2 yields

$$|F(\nu', \mu') - F(\nu, \mu)|^2 \le 8M_D^2 \left(D_h(\nu', \nu) + D_h(\mu', \mu)\right),$$

for all  $\nu', \mu', \nu, \mu \in \mathcal{P}(\mathcal{X})$ . Hence, Assumption 2.1 holds with  $L_F = 8M_D^2$ . Finally, we have

$$|F(\nu,\mu)| \le 2M_D, \quad \left|\frac{\delta F}{\delta \nu}(\nu,\mu,\theta)\right| \le 2M_D, \quad \left|\frac{\delta F}{\delta \mu}(\nu,\mu,w)\right| \le 2M_D,$$

for all  $\nu, \mu \in \mathcal{P}(\mathcal{X})$  and all  $\theta \in \theta, w \in w$ . Hence, Assumption 3.5 holds with  $M = C_{\nu} = C_{\mu} = 2M_{D}$ .

# E EXAMPLE: ADVERSARIAL TRAINING OF MEAN-FIELD NEURAL NETWORKS

Let  $\mathcal{Y} \subset \mathbb{R}$  and  $\mathcal{Z} \subset \mathbb{R}^{d-1}$  be compact with  $\hat{\mu} \in \mathcal{P}(\mathcal{Y} \times \mathcal{Z})$  representing the training data  $(y,z) \in \mathcal{Y} \times \mathcal{Z}$ . Let  $(w,b) \in \mathbb{R}^{d-1} \times \mathbb{R}$  be the parameters of the neural network and let  $\varphi : \mathbb{R} \to \mathbb{R}$  be a bounded, continuous, non-constant activation function. For  $x \coloneqq (w,b) \in \mathbb{R}^d$  and  $z \in \mathbb{R}^{d-1}$ , define the function  $\hat{\varphi}(x,z) \coloneqq \ell(b)\varphi(w \cdot z)$ , where  $\ell : \mathbb{R} \to [-K,K]$  is a clipping function with

clipping threshold K > 0. The training of the two-layer neural network aims to find the optimal set of parameters  $\{x_i\}_{i=1}^N$  which minimize the non-convex  $L^2$ -loss function

$$F_N^0(x_1, ..., x_N) := \frac{1}{2} \int_{\mathcal{Y} \times \mathcal{Z}} \left| y - \frac{1}{N} \sum_{i=1}^N \hat{\varphi}(x_i, z) \right|^2 \hat{\mu}(\mathrm{d}y, \mathrm{d}z). \tag{36}$$

Instead of solving the non-convex minimization problem (36), we lift it to space of probability measures and consider the mean-field optimization problem (see e.g. (Hu et al., 2021, Section 3) and the references therein)

$$\min_{\nu \in \mathcal{P}(\mathbb{R}^d)} F^0(\nu), \quad \text{ with } F^0(\nu) \coloneqq \frac{1}{2} \int_{\mathcal{V} \times \mathcal{Z}} \left| y - \mathbb{E}^{X \sim \nu} [\hat{\varphi}(X,z)] \right|^2 \hat{\mu}(\mathrm{d}y,\mathrm{d}z).$$

To account for potential attacks by an adversary aiming to manipulate the training data  $\hat{\mu}$ , we minimize over the parameter distribution  $\nu$ , considering the "worst-case" perturbation of  $\hat{\mu}$ . This leads to the following mean-field min-max game

$$\min_{\nu \in \mathcal{P}(\mathbb{R}^d)} \max_{\mu \in \mathcal{P}(\mathcal{Y} \times \mathcal{Z})} F^0(\nu, \mu) - \text{TV}^2(\mu, \hat{\mu}), \tag{37}$$

where  $\mathrm{TV}^2$  denotes the squared total variation distance, which represents the cost incurred by the adversary to alter the original training data  $\hat{\mu}$ . The resulting objective  $F(\nu,\mu) := F^0(\nu,\mu) - \mathrm{TV}^2(\mu,\hat{\mu})$  is a non-linear function covered by our general framework. The choice of the incurred cost in (37) is, to an extent, arbitrary, and we focus here on  $\mathrm{TV}^2$  due to its convenience for verifying our assumptions. Alternative cost functions include the Wasserstein distance (Bai et al., 2023; Trillos & Trillos, 2023) and the KL divergence (Si et al., 2023).

**Proposition E.1** (Verification of assumptions for Example E). Let  $\mathcal{Y} \subset \mathbb{R}$  and  $\mathcal{Z} \subset \mathbb{R}^{d-1}$  be compact with  $\hat{\mu} \in \mathcal{P}(\mathcal{Y} \times \mathcal{Z})$ . For  $x := (w,b) \in \mathbb{R}^d$  and  $z \in \mathbb{R}^{d-1}$ , let  $\hat{\varphi}(x,z) := \ell(b)\varphi(w \cdot z)$ , where  $\ell : \mathbb{R} \to [-K,K]$  is a clipping function with clipping threshold K > 0 and  $\varphi : \mathbb{R} \to \mathbb{R}$  is a bounded, continuous, non-constant function. Then Assumptions 1.5, 1.6, 2.1, and 3.5 are satisfied by the objective

$$F(\nu,\mu) = \frac{1}{2} \int_{\mathcal{V} \times \mathcal{Z}} \left| y - \mathbb{E}^{X \sim \nu} [\hat{\varphi}(X,z)] \right|^2 \mu(\mathrm{d}y,\mathrm{d}z) - \mathrm{TV}^2(\mu,\hat{\mu}).$$

*Proof.* Observe that by linearity of the expectation in  $\nu$  and convexity of  $|\cdot|^2$ , the function

$$F^{0}(\nu,\mu) = \frac{1}{2} \int_{\mathcal{V} \times \mathcal{Z}} \left| y - \mathbb{E}^{X \sim \nu} [\hat{\varphi}(X,z)] \right|^{2} \mu(\mathrm{d}y,\mathrm{d}z)$$

satisfies the flat-convexity condition

$$F^{0}((1-\varepsilon)\nu + \varepsilon\nu', \mu) \le (1-\varepsilon)F^{0}(\nu, \mu) + \varepsilon F^{0}(\nu', \mu),$$

for any  $\nu, \nu' \in \mathcal{P}(\mathbb{R}^d)$ ,  $\mu \in \mathcal{P}(\mathcal{Y} \times \mathcal{Z})$  and any  $\varepsilon \in [0, 1]$ . Hence, by (Hu et al., 2021, Lemma 4.1),  $\nu \mapsto F(\nu, \mu)$  satisfies  $D_{F(\cdot, \mu)}(\nu', \nu) \geq 0$ . Again, by convexity of  $|\cdot|^2$ , it holds that  $\mathrm{TV}^2$  is convex, that is,

$$TV^2((1-\varepsilon)\mu + \varepsilon\mu', \hat{\mu}) \le (1-\varepsilon)TV^2(\mu, \hat{\mu}) + \varepsilon TV^2(\mu', \hat{\mu}),$$

for any  $\mu, \mu' \in \mathcal{P}(\mathcal{Y} \times \mathcal{Z})$  and any  $\varepsilon \in [0, 1]$ . Also, by linearity of  $F^0$  in  $\mu$ , it follows that F satisfies the flat concavity condition

$$F(\nu, (1-\varepsilon)\mu + \varepsilon\mu') \ge (1-\varepsilon)F(\nu, \mu) + \varepsilon F(\nu, \mu'),$$

for any  $\mu', \mu \in \mathcal{P}(\mathcal{Y} \times \mathcal{Z}), \nu \in \mathcal{P}(\mathbb{R}^d)$  and any  $\varepsilon \in [0, 1]$ . Hence, by (Hu et al., 2021, Lemma 4.1),  $\mu \mapsto F(\nu, \mu)$  satisfies  $D_{F(\nu, \cdot)}(\mu', \mu) \leq 0$ . Therefore, F satisfies Assumption 1.5.

To verify Assumption 1.6, it is enough to show that for all  $\nu', \nu \in \mathcal{P}(\mathbb{R}^d)$ ,  $\mu \in \mathcal{P}(\mathcal{Y} \times \mathcal{Z})$  and all  $x \in \mathbb{R}^d$ ,

$$\left| \frac{\delta F}{\delta \nu}(\nu', \mu, x) - \frac{\delta F}{\delta \nu}(\nu, \mu, x) \right| \le C_F \operatorname{TV}(\nu', \nu)$$

since by Definition F.1, this implies

$$F(\nu',\mu) - F(\nu,\mu) - \int_{\mathbb{R}^d} \frac{\delta F}{\delta \nu} (\nu,\mu,x) (\nu' - \nu) (\mathrm{d}x)$$

$$= \int_0^1 \int_{\mathbb{R}^d} \left( \frac{\delta F}{\delta \nu} (\nu + \varepsilon(\nu' - \nu), \mu, x) - \frac{\delta F}{\delta \nu} (\nu, \mu, x) \right) (\nu' - \nu) (\mathrm{d}x) \mathrm{d}\varepsilon$$

$$\leq 2C_F \int_0^1 \mathrm{TV} (\nu + \varepsilon(\nu' - \nu), \nu) \, \mathrm{TV} (\nu', \nu) \mathrm{d}\varepsilon$$

$$\leq 2C_F \int_0^1 \varepsilon \, \mathrm{TV}^2 (\nu', \nu) \mathrm{d}\varepsilon$$

$$= C_F \, \mathrm{TV}^2 (\nu', \nu) \leq \frac{C_F}{2} D_h(\nu', \nu),$$

where the last inequality follows from Remark 2.2. Thus,  $D_{F(\cdot,\mu)}(\nu',\nu) \leq L_{\nu}D_{h}(\nu',\nu)$  in Assumption 1.6 holds with  $L_{\nu} = \frac{C_{F}}{2}$ . The same argument applies to  $D_{F(\nu,\cdot)}(\mu',\mu) \geq -L_{\mu}D_{h}(\mu',\mu)$  in Assumption 1.6.

Note that

$$\frac{\delta F}{\delta \nu}(\nu,\mu,x) = -\int_{\mathcal{Y}\times\mathcal{Z}} \left(y - \mathbb{E}^{X\sim\nu}[\hat{\varphi}(X,z)]\right) \hat{\varphi}(x,z) \mu(\mathrm{d}y,\mathrm{d}z).$$

Since  $\varphi$  is bounded by  $M_{\varphi} > 0$ , we obtain

$$\left| \frac{\delta F}{\delta \nu}(\nu', \mu, x) - \frac{\delta F}{\delta \nu}(\nu, \mu, x) \right| \leq \int_{\mathcal{Y} \times \mathcal{Z}} \int_{\mathbb{R}^d} |\hat{\varphi}(x, z)| \, |\nu' - \nu| \, (\mathrm{d}x) \, |\hat{\varphi}(x, z)| \, \mu(\mathrm{d}y, \mathrm{d}z)$$
$$\leq 2K^2 M_{\omega}^2 \, \mathrm{TV}(\nu', \nu).$$

Thus,  $L_{\nu} = K^2 M_{\varphi}^2$ .

Let  $r := (y, z) \in \mathbb{R}^d$ , and assume for simplicity that both  $\mu$ ,  $\hat{\mu}$  are absolutely continuous with respect to Lebesgue measure. We claim that

$$\frac{\delta \operatorname{TV}(\cdot, \hat{\mu})}{\delta \mu}(\mu, r) = \frac{1}{2} \operatorname{sign} \left( \mu(r) - \hat{\mu}(r) \right),$$

for  $\mu \neq \hat{\mu}$  a.e. Fix  $\hat{\mu}$ . For any  $\mu'$ , any  $\mu \neq \hat{\mu}$  a.e., and any  $\varepsilon \in (0,1)$ , (Tsybakov, 2008, Lemma 2.1) gives

$$\lim_{\varepsilon \to 0} \frac{1}{\varepsilon} \left( \text{TV}(\mu + \varepsilon(\mu' - \mu), \hat{\mu}) - \text{TV}(\mu, \hat{\mu}) \right)$$

$$= \lim_{\varepsilon \to 0} \frac{1}{2\varepsilon} \int_{\mathbb{R}^d} \left( |\mu(r) - \hat{\mu}(r) + \varepsilon(\mu'(r) - \mu(r))| - |\mu(r) - \hat{\mu}(r)| \right) dr.$$

Since  $|\cdot|$  is differentiable at every  $v \neq 0$  with derivative  $\operatorname{sign}(v)$ , we obtain by dominated convergence

$$\lim_{\varepsilon \to 0} \frac{1}{\varepsilon} \left( \text{TV}(\mu + \varepsilon(\mu' - \mu), \hat{\mu}) - \text{TV}(\mu, \hat{\mu}) \right) = \frac{1}{2} \int_{\mathbb{R}^d} \text{sign} \left( \mu(r) - \hat{\mu}(r) \right) (\mu'(r) - \mu(r)) (dr).$$

To justify dominated convergence, note that for every r, the reverse triangle inequality gives

$$\left| \frac{|\mu(r) - \hat{\mu}(r) + \varepsilon(\mu'(r) - \mu(r))| - |\mu(r) - \hat{\mu}(r)|}{\varepsilon} \right| \le |\mu'(r) - \mu(r)| \in L^1(\mathbb{R}^d).$$

If  $\mu = \hat{\mu}$  a.e., then the map  $\mathbb{R} \ni v \mapsto |v|$  is not differentiable at v = 0 but its subdifferential is the interval [-1, 1]. Hence, the subdifferential of TV at such measures is the interval  $[-\frac{1}{2}, \frac{1}{2}]$ .

Finally, by the chain rule,

$$\frac{\delta \operatorname{TV}^{2}(\cdot, \hat{\mu})}{\delta \mu}(\mu, r) = 2 \operatorname{TV}(\mu, \hat{\mu}) \frac{\delta \operatorname{TV}(\cdot, \hat{\mu})}{\delta \mu}(\mu, r),$$

and we immediately see that  $\frac{\delta \text{ TV}^2(\cdot,\hat{\mu})}{\delta \mu}(\mu,r) = 0$  if  $\mu = \hat{\mu}$  a.e.. Hence, combining both cases,

$$\frac{\delta \operatorname{TV}^2(\cdot,\hat{\mu})}{\delta \mu}(\mu,r) = \begin{cases} \operatorname{TV}(\mu,\hat{\mu}) \operatorname{sign} \left(\mu(r) - \hat{\mu}(r)\right), & \mu \neq \hat{\mu} \text{ a.e.,} \\ 0, & \mu = \hat{\mu} \text{ a.e.} \end{cases}$$

Consequently,

$$\frac{\delta F}{\delta \mu}(\nu,\mu,r) = \frac{1}{2} \left| y - \mathbb{E}^{X \sim \nu}[\hat{\varphi}(X,z)] \right|^2 - \mathrm{TV}(\mu,\hat{\mu}) \operatorname{sign}\left(\mu(r) - \hat{\mu}(r)\right).$$

Hence.

$$\left| \frac{\delta F}{\delta \mu}(\nu, \mu', ) - \frac{\delta F}{\delta \mu}(\nu, \mu, r) \right| = \left| \text{TV}(\mu', \hat{\mu}) \operatorname{sign} \left( \mu'(r) - \hat{\mu}(r) \right) - \text{TV}(\mu, \hat{\mu}) \operatorname{sign} \left( \mu(r) - \hat{\mu}(r) \right) \right|.$$
(38)

If  $\operatorname{sign}(\mu'(r) - \hat{\mu}(r)) = \operatorname{sign}(\mu(r) - \hat{\mu}(r)) > 0$  a.e. or both are < 0 a.e., then (38) becomes

$$\left| \frac{\delta F}{\delta \mu}(\nu, \mu', r) - \frac{\delta F}{\delta \mu}(\nu, \mu, r) \right| = |\text{TV}(\mu', \hat{\mu}) - \text{TV}(\mu, \hat{\mu})|$$

$$= \frac{1}{2} \left| \int_{\mathbb{R}^d} (\mu'(r) - \hat{\mu}(r)) dr - \int_{\mathbb{R}^d} (\mu(r) - \hat{\mu}(r)) dr \right|$$

$$\leq \text{TV}(\mu', \mu).$$

If  $\operatorname{sign}(\mu'(r) - \hat{\mu}(r)) > 0$  a.e. and  $\operatorname{sign}(\mu(r) - \hat{\mu}(r)) < 0$  a.e., or vice versa, then (38) becomes

$$\begin{split} \left| \frac{\delta F}{\delta \mu}(\nu, \mu', r) - \frac{\delta F}{\delta \mu}(\nu, \mu, r) \right| &= \mathrm{TV}(\mu', \hat{\mu}) + \mathrm{TV}(\mu, \hat{\mu}) \\ &= \frac{1}{2} \int_{\mathbb{R}^d} |\mu'(r) - \hat{\mu}(r)| \, \mathrm{d}r + \frac{1}{2} \int_{\mathbb{R}^d} |\mu(r) - \hat{\mu}(r)| \, \mathrm{d}r \\ &= \frac{1}{2} \int_{\mathbb{R}^d} (\mu'(r) - \hat{\mu}(r)) \, \mathrm{d}r + \frac{1}{2} \int_{\mathbb{R}^d} (\hat{\mu}(r) - \mu(r)) \, \mathrm{d}r \\ &< \mathrm{TV}(\mu', \mu). \end{split}$$

Thus,  $L_{\mu} = \frac{1}{2}$ .

To verify Assumption 2.1 and 3.5, note that

$$\frac{\delta F}{\delta \nu}(\nu, \mu, x) \left| \le \int_{\mathcal{Y} \times \mathcal{Z}} \left| y - \mathbb{E}^{X \sim \nu} [\hat{\varphi}(X, z)] \right| \left| \hat{\varphi}(x, z) \right| \mu(\mathrm{d}y, \mathrm{d}z) \\
\le K M_{\varphi} \left( \mu_{\mathcal{Y}} + K M_{\varphi} \right) := C_{\nu},$$

where

$$\mu_{\mathcal{Y}} := \int_{\mathcal{V} \times \mathcal{Z}} |y| \mu(\mathrm{d}y, \mathrm{d}z) < \infty$$

since  $\mathcal{Y} \times \mathcal{Z}$  is compact.

Similarly,

$$\begin{split} \left| \frac{\delta F}{\delta \mu}(\nu, \mu, r) \right| &= \left| \frac{1}{2} \left| y - \mathbb{E}^{X \sim \nu} [\hat{\varphi}(X, z)] \right|^2 - \text{TV}(\mu, \hat{\mu}) \operatorname{sign} \left( \mu(r) - \hat{\mu}(r) \right) \right| \\ &\leq 1 + \frac{1}{2} \left( \operatorname{diam}(\mathcal{Y}) + K M_{\varphi} \right)^2 \coloneqq C_{\mu}, \end{split}$$

since  $\mathcal{Y}$  is compact and  $TV(\mu, \hat{\mu}) \leq 1$ .

For  $\nu', \nu \in \mathcal{P}(\mathbb{R}^d)$  and  $\mu', \mu \in \mathcal{P}(\mathcal{Y} \times \mathcal{Z})$ ,  $|F(\nu', \mu') - F(\nu, \mu)| = |F(\nu', \mu') - F(\nu, \mu') + F(\nu, \mu') - F(\nu, \mu)|$   $\leq |F(\nu', \mu') - F(\nu, \mu')| + |F(\nu, \mu') - F(\nu, \mu)|$   $= \left| \int_0^1 \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} \left( \nu + \varepsilon(\nu' - \nu), \mu', x \right) (\nu' - \nu) (\mathrm{d}x) \mathrm{d}\varepsilon \right|$   $+ \left| \int_0^1 \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} \left( \nu, \mu + \varepsilon(\mu' - \mu), y \right) (\mu' - \mu) (\mathrm{d}y) \mathrm{d}\varepsilon \right|$   $\leq 2C_{\nu} \operatorname{TV}(\nu', \nu) + 2C_{\mu} \operatorname{TV}(\mu', \mu).$ 

Squaring both sides and using Remark 2.2 gives that Assumption 2.1 holds with  $L_F = 8 \max\{C_{\nu}^2, C_{\mu}^2\}$ 

For Assumption 3.5, observe that

$$|F(\nu,\mu)| \leq \frac{1}{2} + \frac{1}{2} \left( \operatorname{diam}(\mathcal{Y}) + KM_{\varphi} \right)^2 \coloneqq M,$$

since  $\mathcal{Y}$  is compact and  $TV(\mu, \hat{\mu}) \leq 1$ .

# F DIFFERENTIABILITY ON THE PRIMAL SPACE

In this section, following (Carmona & Delarue, 2018, Definition 5.43) and (Santambrogio, 2015, Definition 7.12), we introduce the notion of differentiability on the space of probability measure that we utilize throughout the paper.

**Definition F.1.** For any  $\mathcal{X} \subset \mathbb{R}^d$ , let  $\mathcal{K} \subseteq \mathcal{P}(\mathcal{X})$  be convex. A function  $F: \mathcal{P}(\mathcal{X}) \to \mathbb{R}$  admits first-order flat derivative on  $\mathcal{K}$ , if there exists a measurable function  $\frac{\delta F}{\delta \nu}: \mathcal{K} \times \mathcal{X} \to \mathbb{R}$  such that, for any  $\nu, \nu' \in \mathcal{K}$ , there exists C > 0 such that, for all  $x \in \mathcal{X}$ , we have  $\left|\frac{\delta F}{\delta \nu}(\nu, x)\right| \leq C$ , and it holds that

$$\lim_{\varepsilon \to 0} \frac{F(\nu + \varepsilon(\nu' - \nu)) - F(\nu)}{\varepsilon} = \int_{\mathcal{V}} \frac{\delta F}{\delta \nu}(\nu, x) (\nu' - \nu) (\mathrm{d}x). \tag{39}$$

The functional  $\frac{\delta F}{\delta \nu}$  is called the flat derivative of F on K. We note that  $\frac{\delta F}{\delta \nu}$  exists up to an additive constant, and thus we make the normalizing convention  $\int_{\mathcal{X}} \frac{\delta F}{\delta \nu}(\nu,x) \nu(\mathrm{d}x) = 0$ .

If, for any fixed  $x \in \mathcal{X}$ , the map  $\nu \mapsto \frac{\delta F}{\delta \nu}(\nu, x)$  satisfies Definition F.1, we say that F admits a second-order flat derivative denoted by  $\frac{\delta^2 F}{\delta \nu^2}$ . Consequently, by Definition F.1, there exists a measurable functional  $\frac{\delta^2 F}{\delta \nu^2} : \mathcal{K} \times \mathcal{X} \times \mathcal{X} \to \mathbb{R}$  such that

$$\lim_{\varepsilon \to 0} \frac{1}{\varepsilon} \left( \frac{\delta F}{\delta \nu} (\nu + \varepsilon (\nu' - \nu), x) - \frac{\delta F}{\delta \nu} (\nu, x) \right) = \int_{\mathcal{X}} \frac{\delta^2 F}{\delta \nu^2} (\nu, x, x') (\nu' - \nu) (\mathrm{d}x'). \tag{40}$$

#### G DIFFERENTIABILITY ON THE DUAL SPACE

In this section, we start by recalling the notions of Fréchet and Gâteaux derivative for functions  $H: B_b(\mathcal{X}) \to \mathfrak{X}$ , where  $(B_b(\mathcal{X}), \|\cdot\|_{\infty})$  is the Banach space of real-valued bounded measurable functions on  $\mathcal{X} \subset \mathbb{R}^d$  and  $(\mathfrak{X}, \|\cdot\|_{\mathfrak{X}})$  is a normed vector space; see e.g. Chapters 7, 1, 3 in (Aliprantis & Border, 2007; Ambrosetti & Prodi, 1995; Ortega & Rheinboldt, 1970), respectively. Based on these notions of differentiablity, we will introduce the notions of first and second variation for functions H.

# G.1 Preliminaries on Fréchet and Gâteaux derivatives

For  $\mathcal{X} \subset \mathbb{R}^d$ , let  $\mathcal{L}(B_b(\mathcal{X}), \mathfrak{X})$  and  $\mathcal{L}(B_b(\mathcal{X}))$  denote the space of continuous linear maps from  $B_b(\mathcal{X})$  to  $\mathfrak{X}$ , and from  $B_b(\mathcal{X})$  to itself, respectively.

**Definition G.1** (Fréchet differentiability). Let  $\mathcal{U} \subset B_b(\mathcal{X})$  be open. Given  $f \in \mathcal{U}$ , the function  $H: \mathcal{U} \to \mathfrak{X}$  is Fréchet differentiable at f if there exists  $T \in \mathcal{L}(B_b(\mathcal{X}), \mathfrak{X})$  such that, for all  $g \in B_b(\mathcal{X})$ ,

$$\lim_{\left\Vert g\right\Vert _{\infty}\rightarrow0}\frac{\left\Vert H\left( f+g\right) -H(f)-T\left( g\right) \right\Vert _{\mathfrak{X}}}{\left\Vert g\right\Vert _{\infty}}=0.$$

If it exists, the map T is unique, we write  $T = \nabla_{\mathcal{F}} H(f)$ , and call  $\nabla_{\mathcal{F}} H(f)$  the Fréchet derivative of H at f. If H is Fréchet differentiable at every  $f \in \mathcal{U}$ , then we say that H is Fréchet differentiable on  $\mathcal{U}$ .

**Example G.2** (Convex conjugate of entropy). If h is the entropy, then a straightforward calculation directly from Definition 3.1 shows that its dual  $h^*$  is given by

$$h^*(f) = \log \left( \int_{\mathcal{X}} e^{f(z)} dz \right).$$

**Example G.3** (Fréchet derivative of entropy). From Definition G.1 and following the argument from (Kerimkulov et al., 2023, Proposition 3.9), we can show that  $h^*$  is Fréchet differentiable on  $B_b(\mathcal{X})$  with Fréchet derivative given by

$$\nabla_{\mathcal{F}}h^*(f)(g) = \int_{\mathcal{X}} g(z) \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} dz} dz, \tag{41}$$

for all  $g \in B_b(\mathcal{X})$ .

**Definition G.4** (Gâteaux differentiability). Let  $\mathcal{U} \subset B_b(\mathcal{X})$  be open. Given  $f \in \mathcal{U}$ , the function  $H: \mathcal{U} \to \mathfrak{X}$  is Gâteaux differentiable at f if there exists  $T \in \mathcal{L}(B_b(\mathcal{X}), \mathfrak{X})$  such that for any direction  $f' \in B_b(\mathcal{X})$ ,

$$\lim_{\varepsilon \downarrow 0} \frac{H(f + \varepsilon f') - H(f)}{\varepsilon} = T(f').$$

If it exists, the map T is unique, we write  $T = \nabla_{\mathcal{G}} H(f)$ , and call  $\nabla_{\mathcal{G}} H(f)$  the Gâteaux derivative of H at f. If H is Gâteaux differentiable at every  $f \in \mathcal{U}$ , then we say that H is Gâteaux differentiable on  $\mathcal{U}$ .

As observed in Chapter 1, 3 in (Ambrosetti & Prodi, 1995; Ortega & Rheinboldt, 1970), if H is Fréchet differentiable, then it is automatically Gâteaux differentiable and the two derivatives coincide, i.e.,  $\nabla_{\mathcal{F}}H = \nabla_{\mathcal{G}}H$ . Moreover, (Ortega & Rheinboldt, 1970, Proposition 3.1.6) proves that Fréchet differentiability of H at  $f \in \mathcal{U}$  implies that H is continuous at f, whereas in the case of Gâteaux differentiability, this does not necessarily hold; see (Ortega & Rheinboldt, 1970, Proposition 3.1.4).

Following the discussions in (Aliprantis & Border, 2007; Ambrosetti & Prodi, 1995; Ortega & Rheinboldt, 1970), it is possible to extend Definition G.1 to higher-order Fréchet derivatives.

**Definition G.5** (Second-order Fréchet differentiability). Let  $\mathcal{U} \subset B_b(\mathcal{X})$  be open and let  $f \in \mathcal{U}$ . Suppose that  $H: \mathcal{U} \to \mathfrak{X}$  is Fréchet differentiable (cf. Definition G.1) at f, and admits Fréchet derivative  $\nabla_{\mathcal{F}} H(f)$ . Then  $\nabla_{\mathcal{F}} H(f)$  is Fréchet differentiable at f, if there exists  $T \in \mathcal{L}(B_b(\mathcal{X}), \mathcal{L}(B_b(\mathcal{X}), \mathfrak{X}))$  such that for all  $f', f'' \in B_b(\mathcal{X})$ ,

$$\lim_{\left\|f''\right\|_{\infty}\to 0}\frac{\left\|\nabla_{\mathcal{F}}H\left(f+f''\right)\left(f'\right)-\nabla_{\mathcal{F}}H(f)(f')-T\left(f''\right)\left(f'\right)\right\|_{\mathfrak{X}}}{\left\|f''\right\|_{\infty}}=0.$$

If it exists, the map T is unique, we write  $T = \nabla^2_{\mathcal{F}} H(f)$ , and call  $\nabla^2_{\mathcal{F}} H(f)$  the second Fréchet derivative of H at f.

**Example G.6** (Second order Fréchet derivative of entropy). If h is the entropy, using (41) and following the argument from (Kerimkulov et al., 2023, Proposition 3.6), we can show that  $\nabla_{\mathcal{F}} h^*(f)$  is Fréchet differentiable on  $B_b(\mathcal{X})$  with Fréchet derivative given by

$$\begin{split} &\nabla_{\mathcal{F}}^2 h^*(f)(f')(g) = \int_{\mathcal{X}} g(x) \left( f'(x) - \int_{\mathcal{X}} f'(z) \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}z \right) \frac{e^{f(x)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}x \\ &= \int_{\mathcal{X}} \left( g(x) - \int_{\mathcal{X}} g(z) \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}z + \int_{\mathcal{X}} g(z) \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}z \right) \times \\ &\times \left( f'(x) - \int_{\mathcal{X}} f'(z) \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}z \right) \frac{e^{f(x)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}x \\ &= \int_{\mathcal{X}} \left( g(x) - \int_{\mathcal{X}} g(z) \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}z \right) \left( f'(x) - \int_{\mathcal{X}} f'(z) \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}z \right) \frac{e^{f(x)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}x, \end{split}$$

for all  $g \in B_b(\mathcal{X})$ , where the last line used the fact that

$$\int_{\mathcal{X}} \int_{\mathcal{X}} g(z) \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} dz} \left( f'(x) - \int_{\mathcal{X}} f'(z) \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} dz} dz \right) \frac{e^{f(x)}}{\int_{\mathcal{X}} e^{f(z)} dz} dz dx = 0.$$

The motivation behind working with Fréchet instead of Gâteaux differentiability is that the higherorder derivatives in the case of the former could be identified with continuous symmetric multilinear maps. As proved in Section 3 of Chapter 1 from (Ambrosetti & Prodi, 1995), the space  $\mathcal{L}(B_b(\mathcal{X}), \mathcal{L}(B_b(\mathcal{X}), \mathfrak{X}))$  is isometrically isomorphic to  $\mathcal{L}_2(B_b(\mathcal{X}), \mathfrak{X})$ , i.e., the space of continuous bilinear maps from  $B_b(\mathcal{X}) \times B_b(\mathcal{X})$  to  $\mathfrak{X}$ , and therefore, we could naturally view the second-order Fréchet derivative of H, if it exists, as a continuous bilinear map.

Furthermore, due to (Ambrosetti & Prodi, 1995, Theorem 3.5), we have that the second-order Fréchet derivative is always symmetric. On the contrary, the second-order Gâteaux derivative is not necessarily symmetric as noted on page 78 in (Ortega & Rheinboldt, 1970).

**Remark G.7.** If we replace  $B_b(\mathcal{X})$  with  $\mathbb{R}^d$  and  $\mathfrak{X}$  with  $\mathbb{R}$ , then the first and second-order Fréchet derivatives are precisely the gradient and Hessian matrix of H at f.

# G.2 FIRST AND SECOND VARIATIONS

Following Chapter 2 from (Abraham et al., 2012), we introduce the notions of first and second variation for Fréchet differentiable functions H, relative to the duality pairing (5).

**Definition G.8** (First variation of H). Let  $H: B_b(\mathcal{X}) \to \mathfrak{X}$  be Fréchet differentiable at  $f \in B_b(\mathcal{X})$ . If it exists, the first variation of H at f is the unique continuous map  $B_b(\mathcal{X}) \ni f \mapsto \frac{\delta H}{\delta f}(f) \in \mathcal{P}(\mathcal{X})$  such that, for all  $g \in B_b(\mathcal{X})$ ,

$$\left\langle g, \frac{\delta H}{\delta f}(f) \right\rangle \coloneqq \nabla_{\mathcal{F}} H(f)(g).$$

**Example G.9** (First variation of the dual of entropy). From Example G.2, we observe that the first variation  $\frac{\delta h^*}{\delta f}: B_b(\mathcal{X}) \to \mathcal{P}(\mathcal{X})$  of  $h^*$  is given by

$$\frac{\delta h^*}{\delta f}(f)(\mathrm{d}z) = \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}z.$$

Assuming that  $H: B_b(\mathcal{X}) \to \mathbb{R}$  is Fréchet differentiable at  $f \in B_b(\mathcal{X})$  with Fréchet derivative  $\nabla_{\mathcal{F}} H(f)$ , then it is Gâteaux differentiable (cf. Definition G.4) with the same derivative, and therefore the first variation of H at f can be characterized as

$$\left\langle g, \frac{\delta H}{\delta f}(f) \right\rangle = \lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left( H \left( f + \varepsilon g \right) - H \left( f \right) \right),$$
 (42)

for all  $q \in B_b(\mathcal{X})$ .

Let  $f, g \in B_b(\mathcal{X})$ . For any  $\lambda \in [0, 1]$ , set  $f^{\lambda} := f + \lambda g$ . Then since  $f^{\lambda} \in B_b(\mathcal{X})$ , for all  $\lambda \in [0, 1]$ , it follows by (42) that

$$\lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left( H \left( f^{\lambda} + \varepsilon g \right) - H \left( f^{\lambda} \right) \right) = \left\langle g, \frac{\delta H}{\delta f} \left( f^{\lambda} \right) \right\rangle.$$

Since  $f^{\lambda} + \varepsilon g = f^{\lambda+\varepsilon}$ , it follows by the fundamental theorem of calculus that

$$H(f+g) - H(f) = \int_0^1 \lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left( H\left(f^{\lambda + \varepsilon}\right) - H\left(f^{\lambda}\right) \right) d\lambda = \int_0^1 \left\langle g, \frac{\delta H}{\delta f} \left(f^{\lambda}\right) \right\rangle d\lambda. \tag{43}$$

With the definition of first variation at hand, we can introduce necessary and sufficient conditions for H to have an extremum at  $f \in B_b(\mathcal{X})$ .

**Lemma G.10** (Necessary first-order condition on  $B_b(\mathcal{X})$ ). Let  $\mathfrak{X} = \mathbb{R}$ . Suppose that  $H: B_b(\mathcal{X}) \to \mathbb{R}$  admits first variation at f. If H has an extremum at f, then it holds that

$$\frac{\delta H}{\delta f}(f) = 0.$$

*Proof.* For a proof, see (Abraham et al., 2012, Proposition 2.4.22).

**Lemma G.11** (Sufficient first-order condition on  $B_b(\mathcal{X})$ ). Let  $\mathcal{U} \subset B_b(\mathcal{X})$  be non-empty and convex. Suppose that  $H: \mathcal{U} \to \mathbb{R}$  admits first variation on  $\mathcal{U}$  and is convex in the sense that, for all  $\lambda \in [0,1]$ , and all  $f,g \in \mathcal{U}$ , it holds that  $H((1-\lambda)f + \lambda g) \leq (1-\lambda)H(f) + \lambda H(g)$ . If  $\frac{\delta H}{\delta f}(f^*) = 0$ , for some  $f^* \in \mathcal{U}$ , then  $f^*$  is a global minimum of H.

**Remark G.12.** An analogous result can be identically proved for concave functions and global maxima, so we will give the proof only for the convex case.

*Proof.* Since H is convex and admits first variation, following the argument in (Hu et al., 2021, Lemma 4.1), it can be shown that for any  $f, g \in \mathcal{U}$ 

$$H(g) \ge H(f) + \left\langle g - f, \frac{\delta H}{\delta f}(f) \right\rangle.$$

For  $f = f^*$  and using the assumption that  $\frac{\delta H}{\delta f}(f^*) = 0$ , we get

$$H(g) \ge H(f^*),$$

for all  $g \in \mathcal{U}$ , i.e.  $f^*$  is a global minimum.

**Definition G.13** (Second variation of H). Let  $H: B_b(\mathcal{X}) \to \mathfrak{X}$  be twice Fréchet differentiable at  $f \in B_b(\mathcal{X})$ . If it exists, the second variation of H at f is the unique element  $\frac{\delta^2 H}{\delta f^2}(f) \in \mathcal{L}_2(B_b(\mathcal{X}), \mathfrak{X})$  such that, for all  $q, q' \in B_b(\mathcal{X})$ ,

$$\int_{\mathcal{X}\times\mathcal{X}} g(x) \frac{\delta^2 H}{\delta f^2}(f)(f)(\mathrm{d} y \otimes \mathrm{d} x) g'(y) := \nabla_{\mathcal{F}}^2 H(f)(g)(g'),$$

where  $\frac{\delta^2 H}{\delta f^2}(f)(f)(\mathrm{d}y\otimes\mathrm{d}x)\coloneqq\frac{\delta^2 H}{\delta f^2}(f)(\mathrm{d}x)(f)(\mathrm{d}y).$ 

**Example G.14** (Second variation of the dual of entropy). From Example G.6, we observe that the second variation  $\frac{\delta^2 h^*}{\delta t^2}$ :  $B_b(\mathcal{X}) \to \mathcal{L}(B_b(\mathcal{X}), \mathcal{M}(\mathcal{X}))$  of  $h^*$  is given by

$$\frac{\delta^2 h^*}{\delta f^2}(f)(g)(\mathrm{d}x) = \left(g(x) - \int_{\mathcal{X}} g(z) \frac{e^{f(z)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}z\right) \frac{e^{f(x)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}x. \tag{44}$$

Assume that  $H: B_b(\mathcal{X}) \to \mathfrak{X}$  is twice Fréchet differentiable at  $f \in B_b(\mathcal{X})$ . Then its first-order Fréchet derivative  $\nabla_{\mathcal{F}} H(f)$  is Fréchet differentiable at f, and thus it is Gâteaux differentiable (cf. Definition G.4) with the same second-order derivative. Hence, using Definition G.8, the second variation of H at f can be characterized in terms of the first variation as

$$\int_{\mathcal{X}\times\mathcal{X}} g(x) \frac{\delta^2 H}{\delta f^2}(f)(f)(\mathrm{d}y\otimes\mathrm{d}x)g'(y) = \lim_{\varepsilon\downarrow 0} \frac{1}{\varepsilon} \left\langle g, \left(\frac{\delta H}{\delta f}(f+\varepsilon g') - \frac{\delta H}{\delta f}(f)\right) \right\rangle, \tag{45}$$

for all  $g, g' \in B_b(\mathcal{X})$ .

Let  $f, g, g' \in B_b(\mathcal{X})$ . For any  $\lambda \in [0, 1]$ , set  $f^{\lambda} := f + \lambda g'$ . Then since  $f^{\lambda} \in B_b(\mathcal{X})$ , for all  $\lambda \in [0, 1]$ , it follows by (45) that

$$\lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left\langle g, \left( \frac{\delta H}{\delta f} \left( f^{\lambda} + \varepsilon g' \right) - \frac{\delta H}{\delta f} \left( f^{\lambda} \right) \right) \right\rangle = \int_{\mathcal{X} \times \mathcal{X}} g(x) \frac{\delta^{2} H}{\delta f^{2}} \left( f^{\lambda} \right) \left( dy \otimes dx \right) g'(y).$$

Since  $f^{\lambda} + \varepsilon g' = f^{\lambda + \varepsilon}$ , it follows that

$$\left\langle g, \left( \frac{\delta H}{\delta f} \left( f + g' \right) - \frac{\delta H}{\delta f} \left( f \right) \right) \right\rangle = \left\langle g, \int_{0}^{1} \lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left( \frac{\delta H}{\delta f} \left( f^{\lambda + \varepsilon} \right) - \frac{\delta H}{\delta f} \left( f^{\lambda} \right) \right) d\lambda \right\rangle$$

$$= \int_{0}^{1} \lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} \left\langle g, \left( \frac{\delta H}{\delta f} \left( f^{\lambda + \varepsilon} \right) - \frac{\delta H}{\delta f} \left( f^{\lambda} \right) \right) \right\rangle d\lambda \qquad (46)$$

$$= \int_{0}^{1} \int_{\mathcal{X} \times \mathcal{X}} g(x) \frac{\delta^{2} H}{\delta f^{2}} \left( f^{\lambda} \right) \left( f^{\lambda} \right) (dy \otimes dx) g'(y) d\lambda,$$

where the first equality follows from the fundamental theorem of calculus and the second equality from Fubini's theorem and the dominated convergence theorem.

**Proposition G.15** (Verification of Assumption 3.4 for entropy). Suppose that  $\mathcal{X} \subset \mathbb{R}^d$  is bounded. Let  $f, g \in B_b(\mathcal{X})$  and denote  $\varphi(f)(\mathrm{d}x) \coloneqq \frac{e^{f(x)}}{\int_{\mathcal{X}} e^{f(z)} \mathrm{d}z} \mathrm{d}x \in \mathcal{P}(\mathcal{X})$ . Then, for h being the entropy, its second variation (44) satisfies Assumption 3.4.

*Proof.* Note that the second variation (44) can be written as

$$\frac{\delta^2 h^*}{\delta f^2}(f)(g)(\mathrm{d}x) = \left(g(x) - \int_{\mathcal{X}} g(z)\varphi(f)(\mathrm{d}z)\right)\varphi(f)(\mathrm{d}x).$$

Since  $\varphi(f)$  is absolutely continuous with respect to Lebesgue measure on  $\mathcal{X}$ , it follows that

$$\begin{split} & \left\| \frac{\delta^2 h^*}{\delta f^2}(f')(g') - \frac{\delta^2 h^*}{\delta f^2}(f)(g) \right\|_{\text{TV}} = \int_{\mathcal{X}} \left| \frac{\delta^2 h^*}{\delta f^2}(f')(g') - \frac{\delta^2 h^*}{\delta f^2}(f)(g) \right| (\mathrm{d}x) \\ &= \int_{\mathcal{X}} \left| \left( g'(x) - \int_{\mathcal{X}} g'(z) \varphi(f')(z) \mathrm{d}z \right) \varphi(f')(x) - \left( g(x) - \int_{\mathcal{X}} g(z) \varphi(f)(z) \mathrm{d}z \right) \varphi(f)(x) \right| \mathrm{d}x \\ &= \int_{\mathcal{X}} \left| g'(x) \varphi(f')(x) - \varphi(f')(x) \int_{\mathcal{X}} g'(z) \varphi(f')(z) \mathrm{d}z - g(x) \varphi(f)(x) + \varphi(f)(x) \int_{\mathcal{X}} g(z) \varphi(f)(z) \mathrm{d}z \right| \mathrm{d}x \\ &\leq I_1 + I_2, \end{split}$$

where

$$I_{1} = \int_{\mathcal{X}} |g'(x)\varphi(f')(x) - g(x)\varphi(f)(x)| \, \mathrm{d}x,$$

$$I_{2} = \int_{\mathcal{X}} \left| \varphi(f)(x) \int_{\mathcal{X}} g(z)\varphi(f)(z) \, \mathrm{d}z - \varphi(f')(x) \int_{\mathcal{X}} g'(z)\varphi(f')(z) \, \mathrm{d}z \right| \, \mathrm{d}x,$$

and the last inequality follows from triangle inequality. For  $I_1$ , we observe that

$$I_{1} = \int_{\mathcal{X}} |g'(x)\varphi(f')(x) - g'(x)\varphi(f)(x) + g'(x)\varphi(f)(x) - g(x)\varphi(f)(x)| dx$$

$$\leq \int_{\mathcal{X}} |g'(x)| |\varphi(f')(x) - \varphi(f)(x)| dx + \int_{\mathcal{X}} |\varphi(f)(x)| |g'(x) - g(x)| dx.$$

Since  $f, g' \in B_b(\mathcal{X})$ , there exist  $C_{g'}, C_f > 0$  such that  $|g'(x)| \leq C_{g'}$  and  $|\varphi(f)(x)| \leq C_f$ , for all  $x \in \mathcal{X}$ . Since f, f' are bounded on  $\mathcal{X}$ , following the argument in (Lascu et al., 2025, Lemma A.2), we deduce that  $f \mapsto \varphi(f)$  is Lipschitz, i.e., there exists  $L_{\varphi} > 0$  such that, for all  $x \in \mathcal{X}$ ,

$$|\varphi(f')(x) - \varphi(f)(x)| \le L_{\varphi} |f'(x) - f(x)|.$$

Hence,  $I_1$  becomes

$$I_{1} \leq C_{g'}L_{\varphi} \int_{\mathcal{X}} |f'(x) - f(x)| \, \mathrm{d}x + C_{f} \int_{\mathcal{X}} |g'(x) - g(x)| \, \mathrm{d}x$$

$$\leq \max\{C_{g'}L_{\varphi}, C_{f}\} \int_{\mathcal{X}} (|f'(x) - f(x)| + |g'(x) - g(x)|) \, \mathrm{d}x$$

$$\leq \max\{C_{g'}L_{\varphi}, C_{f}\} |\mathcal{X}| \left( ||f' - f||_{\infty} + ||g' - g||_{\infty} \right).$$

Similarly, for  $I_2$ , we have that

$$\begin{split} I_2 &= \int_{\mathcal{X}} \left| \varphi(f)(x) \int_{\mathcal{X}} g(z) \varphi(f)(z) \mathrm{d}z - \varphi(f)(x) \int_{\mathcal{X}} g'(z) \varphi(f')(z) \mathrm{d}z \right. \\ &+ \varphi(f)(x) \int_{\mathcal{X}} g'(z) \varphi(f')(z) \mathrm{d}z - \varphi(f')(x) \int_{\mathcal{X}} g'(z) \varphi(f')(z) \mathrm{d}z \right| \mathrm{d}x \\ &\leq \int_{\mathcal{X}} \left| \varphi(f)(x) \right| \int_{\mathcal{X}} \left| g(z) \varphi(f)(z) - g'(z) \varphi(f')(z) \right| \mathrm{d}z \mathrm{d}x \\ &+ \int_{\mathcal{X}} \left| \varphi(f)(x) - \varphi(f')(x) \right| \int_{\mathcal{X}} \left| g'(z) \right| \left| \varphi(f')(z) \right| \mathrm{d}z \mathrm{d}x \\ &\leq C_f \int_{\mathcal{X}} \left| g(z) \varphi(f)(z) - g'(z) \varphi(f')(z) \right| \mathrm{d}z + C_{g'} C_{f'} L_{\varphi} \int_{\mathcal{X}} \left| f(x) - f'(x) \right| \mathrm{d}x \\ &\leq C_f \int_{\mathcal{X}} \left| g(z) \varphi(f)(z) - g'(z) \varphi(f')(z) \right| \mathrm{d}z + C_{g'} C_{f'} L_{\varphi} |\mathcal{X}| \|f' - f\|_{\infty}. \end{split}$$

We observe that

$$C_{f} \int_{\mathcal{X}} |g(z)\varphi(f)(z) - g'(z)\varphi(f')(z)| dz dx$$

$$= C_{f} \int_{\mathcal{X}} |g(z)\varphi(f)(z) - g'(z)\varphi(f)(z) + g'(z)\varphi(f)(z) - g'(z)\varphi(f')(z)| dz dx$$

$$\leq C_{f} \int_{\mathcal{X}} |g(z) - g'(z)||\varphi(f)(z)| dz dx + C_{f} \int_{\mathcal{X}} |g'(z)|\varphi(f)(z) - \varphi(f')(z)| dz dx$$

$$\leq |\mathcal{X}|C_{f}^{2}||g' - g||_{\infty} + |\mathcal{X}|C_{f}C_{g'}L_{\varphi}||f' - f||_{\infty}.$$

Hence, we get that

$$I_{1} + I_{2} \leq |\mathcal{X}| \max\{C_{g'}L_{\varphi}, C_{f}\} (\|f' - f\|_{\infty} + \|g' - g\|_{\infty})$$

$$+ |\mathcal{X}|C_{f}^{2}\|g' - g\|_{\infty} + |\mathcal{X}|C_{f}C_{g'}L_{\varphi}\|f' - f\|_{\infty} + |\mathcal{X}|C_{g'}C_{f'}L_{\varphi}\|f' - f\|_{\infty}$$

$$= |\mathcal{X}| \max\{C_{g'}L_{\varphi}, C_{f}\} (\|f' - f\|_{\infty} + \|g' - g\|_{\infty})$$

$$+ |\mathcal{X}|C_{f}^{2}\|g' - g\|_{\infty} + |\mathcal{X}| (C_{f} + C_{f'}) C_{g'}L_{\varphi}\|f' - f\|_{\infty}$$

$$\leq |\mathcal{X}| \left( \max\{C_{g'}L_{\varphi}, C_{f}\} + \max\{C_{f}^{2}, (C_{f} + C_{f'}) C_{g'}L_{\varphi}\} \right) (\|f' - f\|_{\infty} + \|g' - g\|_{\infty}).$$

Setting  $L_{h^*} \coloneqq |\mathcal{X}| \left( \max\{C_{g'}L_{\varphi}, C_f\} + \max\{C_f^2, (C_f + C_{f'}) C_{g'}L_{\varphi}\} \right)$  finishes the verification.

# H TECHNICAL RESULTS ON DUALITY

In this section we state and prove some technical results which are central to the proof technique via dual Bregman divergence that we developed in Subsection 3.

**Proposition H.1.** Let Assumption 1.1 hold. Let  $h^*: B_b(\mathcal{X}) \to \mathbb{R}$  be the convex conjugate of h. Then, the following are equivalent:

- 1. The supremum of  $\mathcal{E} \ni m \mapsto \langle g^*, m \rangle h(m) \in \mathbb{R}$  is attained at  $m = m^*$ ,
- 2.  $g^*(x) \frac{\delta h}{\delta m}(m^*, x) = constant, for all \ x \in \mathcal{X}$  Lebesgue a.e.,
- 3. The supremum of  $B_b(\mathcal{X}) \ni g \mapsto \langle g, m^* \rangle h^*(g) \in \mathbb{R}$  is attained at  $g = g^*$ ,
- 4.  $m^* = \frac{\delta h^*}{\delta g}(g^*)$ .

*Proof.* (1)  $\Longrightarrow$  (2): Suppose that (1) holds. Then the supremum of  $m \mapsto \langle g^*, m \rangle - h(m)$  is attained at the maximizer  $m^* = \arg\max_{m \in \mathcal{E}} \{\langle g^*, m \rangle - h(m)\}$ . Hence, by (Hu et al., 2021, Proposition 2.5),  $m^*$  satisfies the first-order condition

$$g^*(z) - \frac{\delta h}{\delta m}(m^*, z) = \text{constant},$$

for all  $z \in \mathcal{X}$  Lebesgue a.e.

- (2)  $\Longrightarrow$  (1): Suppose that (2) holds. Observe that the map  $m \mapsto \langle g^*, m \rangle h(m)$  is strictly concave due to the strict convexity of h and the linearity of  $m \mapsto \langle g^*, m \rangle$ . Then by the converse of (Hu et al., 2021, Proposition 2.5), it follows that  $m^*$  is the maximizer of the map  $\mathcal{E} \ni m \mapsto \langle g^*, m \rangle h(m) \in \mathbb{R}$ , and so (1) holds.
- (3)  $\Longrightarrow$  (4): Suppose that (3) holds. Then the supremum in  $g \mapsto \langle g, m^* \rangle h^*(g)$  is attained at a maximizer  $g^* \in \arg\max_{g \in B_b(\mathcal{X})} \{\langle g, m^* \rangle h^*(g)\}$ . Hence, by Lemma G.10, it follows that  $g^*$  satisfies the first-order condition

$$m^* = \frac{\delta h}{\delta g}(g^*).$$

(4)  $\Longrightarrow$  (3): Suppose that (4) holds. Observe that  $B_b(\mathcal{X})$  is convex and the map  $g \mapsto \langle g, m^* \rangle - h^*(g)$  is concave due to the convexity of  $h^*$  and the linearity of  $g \mapsto \langle g, m^* \rangle$ . Hence, by Lemma

G.11, it follows that  $g^*$  is a maximizer of the map  $B_b(\mathcal{X}) \ni g \mapsto \langle g, m^* \rangle - h^*(g) \in \mathbb{R}$ , and so (3) holds.

(1)  $\Longrightarrow$  (3): Suppose that (1) holds. Then, by Definition 3.1, we have that  $h^*(g) = \langle g, m^* \rangle - h(m^*)$ , and equivalently  $h(m^*) = \langle g, m^* \rangle - h^*(g)$ . Clearly,  $\mathcal{P}(\mathcal{X})$  is convex and  $(\mathcal{P}(\mathcal{X}), \mathrm{TV})$ , where TV is the total variation distance, is Hausdorff since it is a metric space, hence we can apply the Fenchel-Moreau theorem (Zalinescu, 2002, Theorem 2.3.3) to conclude that  $h^{**} = h$ , i.e.,  $h(m^*) = \sup_{g \in B_b(\mathcal{X})} \{\langle g, m^* \rangle - h^*(g)\}$ . Therefore,  $h(m^*)$  is the supremum of  $g \mapsto \langle g, m^* \rangle - h^*(g)$  attained at  $g = g^*$ .

(3)  $\Longrightarrow$  (1): Suppose (3) holds. Then  $h^{**}(m^*) = \langle g^*, m^* \rangle - h^*(g^*)$ , or equivalently  $h^*(g^*) = \langle g^*, m^* \rangle - h^{**}(m^*)$ . Again, by the Fenchel-Moreau theorem (Zalinescu, 2002, Theorem 2.3.3),  $h^{**}(m) = h(m)$ , for all  $m \in \mathcal{E}$ , and hence  $h^*(g^*) = \langle g^*, m^* \rangle - h(m^*)$ . Hence, by Definition 3.1, the supremum of  $m \mapsto \langle g^*, m \rangle - h(m)$  is realized at  $m = m^*$ .

**Lemma H.2.** Let Assumption 1.1 hold. Let  $h^*: B_b(\mathcal{X}) \to \mathbb{R}$  be the convex conjugate of h. Fix  $f,g \in B_b(\mathcal{X})$  and  $\mu,\mu' \in \mathcal{E}$ . If  $f(z) = \frac{\delta h}{\delta m}(\mu,z)$  and  $g(z) = \frac{\delta h}{\delta m}(\mu',z)$ , for all  $z \in \mathcal{X}$  Lebesgue a.e., up to an additive constant, then

$$D_{h^*}(f,g) = D_h(\mu',\mu).$$

*Proof.* By Definition 3.3, we have that

$$D_{h^*}(f,g) = h^*(f) - h^*(g) - \int_{\mathcal{X}} (f(z) - g(z)) \frac{\delta h^*}{\delta g}(g)(\mathrm{d}z)$$

$$= \langle f, \mu \rangle - h(\mu) - \langle g, \mu' \rangle + h(\mu') - \int_{\mathcal{X}} (f(z) - g(z)) \frac{\delta h^*}{\delta g}(g)(\mathrm{d}z)$$

$$= h(\mu') - h(\mu) + \int_{\mathcal{X}} \frac{\delta h}{\delta m}(\mu, z)\mu(\mathrm{d}z) - \int_{\mathcal{X}} \frac{\delta h}{\delta m}(\mu', z)\mu'(\mathrm{d}z) - \int_{\mathcal{X}} \left(\frac{\delta h}{\delta m}(\mu, z) - \frac{\delta h}{\delta m}(\mu', z)\right) \mu'(\mathrm{d}z)$$

$$= h(\mu') - h(\mu) - \int_{\mathcal{X}} \frac{\delta h}{\delta m}(\mu, z)(\mu' - \mu)(\mathrm{d}z) = D_h(\mu', \mu),$$

where the second and third equalities follow from Lemma H.1 and Corollary 3.2, while the last equality follows from the definition of the Bregman divergence.  $\Box$ 

**Lemma H.3.** Consider (1) and (2). Let Assumption 1.1 hold. Let  $h^*: B_b(\mathcal{X}) \to \mathbb{R}$  be the convex conjugate of h. For each  $n \geq 0$ , fix  $f^n, g^n \in B_b(\mathcal{X})$ ,  $\nu^n \in \mathcal{C}$  and  $\mu^n \in \mathcal{D}$ . If  $f^n = \frac{\delta h}{\delta \nu}(\nu^n, \cdot)$  and  $g^n = \frac{\delta h}{\delta \mu}(\mu^n, \cdot)$ , then, for any  $n \geq 0$ , we have that

$$D_h(\nu^{n+1}, \nu^n) = D_{h^*}(f^n, f^{n+1}), \quad D_h(\nu^n, \nu^{n+1}) = D_{h^*}(f^{n+1}, f^n),$$
  
$$D_h(\mu^{n+1}, \mu^n) = D_{h^*}(g^n, g^{n+1}), \quad D_h(\mu^n, \mu^{n+1}) = D_{h^*}(g^{n+1}, g^n).$$

*Proof.* First, observe that due to Assumption 1.1, the pairs  $(\nu^{n+1}, \mu^{n+1})$  in (1) and (2) are unique. We will only present the proof for (1) since the argument for (2) is identical. The updates in (1) can be equivalently written as

$$\nu^{n+1} = \arg\min_{\nu \in \mathcal{C}} \left\{ \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) (\nu - \nu^{n}) (\mathrm{d}x) + \frac{1}{\tau} D_{h}(\nu, \nu^{n}) \right\}$$

$$= \arg\min_{\nu \in \mathcal{C}} \left\{ \int_{\mathcal{X}} \tau \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) (\nu - \nu^{n}) (\mathrm{d}x) + h(\nu) - h(\nu^{n}) - \int_{\mathcal{X}} \frac{\delta h}{\delta \nu} (\nu^{n}, x) (\nu - \nu^{n}) (\mathrm{d}x) \right\}$$

$$= \arg\min_{\nu \in \mathcal{C}} \left\{ \int_{\mathcal{X}} \left( \tau \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) - \frac{\delta h}{\delta \nu} (\nu^{n}, x) \right) (\nu - \nu^{n}) (\mathrm{d}x) + h(\nu) \right\}$$

$$= \arg\max_{\nu \in \mathcal{C}} \left\{ \int_{\mathcal{X}} \left( \frac{\delta h}{\delta \nu} (\nu^{n}, x) - \tau \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) \right) (\nu - \nu^{n}) (\mathrm{d}x) - h(\nu) \right\}$$

$$= \arg\max_{\nu \in \mathcal{C}} \left\{ \int_{\mathcal{X}} \left( \frac{\delta h}{\delta \nu} (\nu^{n}, x) - \tau \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) \right) \nu (\mathrm{d}x) - h(\nu) \right\},$$

$$(47)$$

and

$$\mu^{n+1} = \underset{\mu \in \mathcal{D}}{\arg \max} \left\{ \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n}, \mu^{n}, y) (\mu - \mu^{n}) (\mathrm{d}y) - \frac{1}{\tau} D_{h}(\mu, \mu^{n}) \right\}$$

$$= \underset{\mu \in \mathcal{D}}{\arg \max} \left\{ \int_{\mathcal{X}} \tau \frac{\delta F}{\delta \mu} (\nu^{n}, \mu^{n}, y) (\mu - \mu^{n}) (\mathrm{d}y) - h(\mu) + h(\mu^{n}) + \int_{\mathcal{X}} \frac{\delta h}{\delta \mu} (\mu^{n}, y) (\mu - \mu^{n}) (\mathrm{d}y) \right\}$$

$$= \underset{\mu \in \mathcal{D}}{\arg \max} \left\{ \int_{\mathcal{X}} \left( \frac{\delta h}{\delta \mu} (\mu^{n}, y) + \tau \frac{\delta F}{\delta \mu} (\nu^{n}, \mu^{n}, y) \right) (\mu - \mu^{n}) (\mathrm{d}y) - h(\mu) \right\}$$

$$= \underset{\mu \in \mathcal{D}}{\arg \max} \left\{ \int_{\mathcal{X}} \left( \frac{\delta h}{\delta \mu} (\mu^{n}, y) + \tau \frac{\delta F}{\delta \mu} (\nu^{n}, \mu^{n}, y) \right) \mu(\mathrm{d}y) - h(\mu) \right\}.$$

$$(48)$$

Using the notation  $f^n = \frac{\delta h}{\delta \nu}(\nu^n, \cdot)$  and  $g^n = \frac{\delta h}{\delta \mu}(\mu^n, \cdot)$ , for each  $n \geq 0$ , the first-order conditions for (1) can be equivalently written as

$$f^{n+1}(x) - f^n(x) = -\tau \frac{\delta F}{\delta \nu}(\nu^n, \mu^n, x),$$
 (49)

$$g^{n+1}(y) - g^n(y) = \tau \frac{\delta F}{\delta \mu}(\nu^n, \mu^n, y), \tag{50}$$

for all  $(x, y) \in \mathcal{X} \times \mathcal{X}$  Lebesgue a.e. Then using (6), (47) becomes

$$\nu^{n+1} = \underset{\nu \in \mathcal{C}}{\arg\max} \left\{ \int_{\mathcal{X}} \left( f^{n}(x) - \tau \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n}, x) \right) \nu(\mathrm{d}x) - h(\nu) \right\}$$
$$= \underset{\nu \in \mathcal{C}}{\arg\max} \left\{ \int_{\mathcal{X}} f^{n+1}(x) \nu(\mathrm{d}x) - h(\nu) \right\} = \frac{\delta h^{*}}{\delta f} (f^{n+1}), \tag{51}$$

for all  $n \ge 0$ . Similarly, from (48), we have that

$$\mu^{n+1} = \frac{\delta h^*}{\delta f}(g^{n+1}),\tag{52}$$

for all  $n \ge 0$ . The conclusion follows directly from Lemma H.2.

# I PROOF OF CONVERGENCE FOR THE MDA IMPLICIT ALGORITHM

In this section, we prove that an implicit Euler discretization of the Fisher-Rao flows studied in (Lascu et al., 2024) yields a linear convergence rate  $\mathcal{O}(1/N)$ , which matches the result in continuous-time under the same assumption of convexity-concavity of F (see (Lascu et al., 2024, Theorem 2.3)). However, a major weakness of this implicit game is that it is not implementable in practice as opposed to (1) and (2).

For a given stepsize  $\tau > 0$ , and fixed initial pair of strategies  $(\nu_0, \mu_0) \in \mathcal{C} \times \mathcal{D}$ , for  $n \geq 0$ , the *implicit* MDA algorithm is defined by

#### **Algorithm 7:** IMPLICIT MDA

**Theorem I.1** (Convergence of the implicit MDA algorithm (7)). Let  $(\nu^*, \mu^*)$  be an MNE of (1) and  $(\nu^0, \mu^0)$  be such that  $\sup_{\nu \in \mathcal{C}} D_h(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}} D_h(\mu, \mu^0) < \infty$ . Let Assumption 1.1, 1.5 and 1.6 hold. Suppose that  $\tau L \leq 1$ , where  $L := \max\{L_{\nu}, L_{\mu}\}$ . Then, we have

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n+1}\right) \leq \frac{1}{N\tau}\left(\sup_{\nu \in \mathcal{C}}D_h(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}}D_h(\mu, \mu^0)\right).$$

*Proof.* Since  $\nu \mapsto \tau \int \frac{\delta F}{\delta \nu}(\nu^n, \mu^{n+1}, x)(\nu - \nu^n)(\mathrm{d}x)$  is convex, applying Lemma C.1 with  $\bar{\nu} = \nu^{n+1}$  and  $\mu = \nu^n$  implies that, for any  $\nu \in \mathcal{C}$ , we have

$$\tau \int \frac{\delta F}{\delta \nu} (\nu^n, \mu^{n+1}, x) (\nu - \nu^n) (\mathrm{d}x) + D_h(\nu, \nu^n) \ge \tau \int \frac{\delta F}{\delta \nu} (\nu^n, \mu^{n+1}, x) (\nu^{n+1} - \nu^n) (\mathrm{d}x) + D_h(\nu^{n+1}, \nu^n) + D_h(\nu, \nu^{n+1}),$$

or, equivalently,

$$-\tau \int \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n+1}, x) (\nu - \nu^{n}) (\mathrm{d}x) - D_{h}(\nu, \nu^{n}) \le -\tau \int \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n+1}, x) (\nu^{n+1} - \nu^{n}) (\mathrm{d}x) - D_{h}(\nu^{n+1}, \nu^{n}) - D_{h}(\nu, \nu^{n+1}).$$
 (53)

Similarly, since  $\mu \mapsto -\tau \int \frac{\delta F}{\delta \mu}(\nu^{n+1}, \mu^n, y)(\mu - \mu^n)(\mathrm{d}y)$  is convex, applying Lemma C.1 with  $\bar{\nu} = \mu^{n+1}$  and  $\mu = \mu^n$  implies that, for any  $\mu \in \mathcal{D}$ , we have

$$\tau \int \frac{\delta F}{\delta \mu} (\nu^{n+1}, \mu^n, y) (\mu - \mu^n) (\mathrm{d}y) - D_h(\mu, \mu^n) \le \tau \int \frac{\delta F}{\delta \mu} (\nu^{n+1}, \mu^n, y) (\mu^{n+1} - \mu^n) (\mathrm{d}y) - D_h(\mu^{n+1}, \mu^n) - D_h(\mu, \mu^{n+1}). \tag{54}$$

Using the convexity of  $\nu \mapsto F(\nu, \mu)$  in (53), with  $\nu = \nu^n$  and  $\mu = \mu^{n+1}$ , we have that

$$F(\nu^{n}, \mu^{n+1}) - F(\nu, \mu^{n+1}) - \frac{1}{\tau} D_{h}(\nu, \nu^{n}) \leq \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n+1}, x) (\nu^{n} - \nu^{n+1}) (\mathrm{d}x) - \frac{1}{\tau} D_{h}(\nu^{n+1}, \nu^{n}) - \frac{1}{\tau} D_{h}(\nu, \nu^{n+1}).$$
 (55)

From  $L_{\nu}$ -relative smoothness and the fact that  $\tau L \leq 1$ , it follows that

$$F(\nu^{n+1}, \mu^{n+1}) \leq F(\nu^{n}, \mu^{n+1}) + \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n+1}, x) (\nu^{n+1} - \nu^{n}) (\mathrm{d}x) + L_{\nu} D_{h}(\nu^{n+1}, \nu^{n})$$

$$\leq F(\nu^{n}, \mu^{n+1}) + \int_{\mathcal{X}} \frac{\delta F}{\delta \nu} (\nu^{n}, \mu^{n+1}, x) (\nu^{n+1} - \nu^{n}) (\mathrm{d}x) + \frac{1}{\tau} D_{h}(\nu^{n+1}, \nu^{n}). \quad (56)$$

Hence, combining (55) with (56), we obtain that

$$F(\nu^{n}, \mu^{n+1}) - F(\nu, \mu^{n+1}) - \frac{1}{\tau} D_{h}(\nu, \nu^{n}) \le F(\nu^{n}, \mu^{n+1}) - F(\nu^{n+1}, \mu^{n+1}) - \frac{1}{\tau} D_{h}(\nu, \nu^{n+1}).$$
(57)

Similarly, using concavity of  $\mu \mapsto F(\nu, \mu)$  in (54), with  $\nu = \nu^{n+1}$  and  $\mu = \mu^n$ , we have that

$$F(\nu^{n+1}, \mu) - F(\nu^{n+1}, \mu^n) - \frac{1}{\tau} D_h(\mu, \mu^n) \le \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n+1}, \mu^n, y) (\mu^{n+1} - \mu^n) (\mathrm{d}y) - \frac{1}{\tau} D_h(\mu^{n+1}, \mu^n) - \frac{1}{\tau} D_h(\mu, \mu^{n+1}).$$
 (58)

From  $L_{\mu}$ -relative smoothness and the fact that  $\tau L \leq 1$ , it follows that

$$F(\nu^{n+1}, \mu^{n+1}) \ge F(\nu^{n+1}, \mu^n) + \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n+1}, \mu^n, y) (\mu^{n+1} - \mu^n) (\mathrm{d}y) - L_{\mu} D_h(\mu^{n+1}, \mu^n)$$

$$\ge F(\nu^{n+1}, \mu^n) + \int_{\mathcal{X}} \frac{\delta F}{\delta \mu} (\nu^{n+1}, \mu^n, y) (\mu^{n+1} - \mu^n) (\mathrm{d}y) - \frac{1}{\tau} D_h(\mu^{n+1}, \mu^n). \tag{59}$$

Hence, combining (58) with (59), we obtain that

$$F(\nu^{n+1},\mu) - F(\nu^{n+1},\mu^n) - \frac{1}{\tau} D_h(\mu,\mu^n) \le F(\nu^{n+1},\mu^{n+1}) - F(\nu^{n+1},\mu^n) - \frac{1}{\tau} D_h(\mu,\mu^{n+1}).$$
(60)

Adding inequalities (57) and (60) implies that

$$F(\nu^{n+1},\mu) - F(\nu,\mu^{n+1}) \le F(\nu^{n+1},\mu^{n+1}) - F(\nu^{n+1},\mu^{n+1})$$
  
 
$$+ \frac{1}{\tau} D_h(\nu,\nu^n) + \frac{1}{\tau} D_h(\mu,\mu^n) - \frac{1}{\tau} D_h(\nu,\nu^{n+1}) - \frac{1}{\tau} D_h(\mu,\mu^{n+1})$$

Summing the previous inequality over n=0,1,...,N-1, bounding the right-hand side from above by its supremum over  $(\nu,\mu)$ , dividing by N, applying Jensen's inequality and taking maximum over  $(\nu,\mu)$  in the left-hand side leads to

$$\operatorname{NI}\left(\frac{1}{N}\sum_{n=0}^{N-1}\nu^{n+1}, \frac{1}{N}\sum_{n=0}^{N-1}\mu^{n+1}\right) \leq \frac{1}{N\tau}\left(\sup_{\nu \in \mathcal{C}}D_h(\nu, \nu^0) + \sup_{\mu \in \mathcal{D}}D_h(\mu, \mu^0)\right),$$

where the last inequality follows since  $D_h(\nu, \nu^N) + D_h(\mu, \mu^N) \ge 0$ , for all  $(\nu, \mu) \in \mathcal{C} \times \mathcal{D}$ .

# J FURTHER RELATED WORKS

Besides the vanilla MDA algorithm, (Hsieh et al., 2019) considers the entropic Mirror Prox algorithm, which requires the computation of an extra gradient at an intermediate point and two projections onto the dual space. Although it is proved in (Hsieh et al., 2019) that the Mirror Prox algorithm achieves  $\mathcal{O}\left(N^{-1}\right)$  convergence rate for deterministic gradients, it is also outlined that for stochastic gradients (which one has typically access to in practice) Mirror Prox and simultaneous MDA achieve the same rate  $\mathcal{O}\left(N^{-1/2}\right)$ .

Another approach based on reproducing kernel Hilbert spaces (RKHS) is developed in (Dvurechensky & Zhu, 2024) and achieves the same convergence rates  $\mathcal{O}\left(N^{-1}\right)$  and  $\mathcal{O}\left(N^{-1/2}\right)$  for the deterministic and stochastic Mirror Prox algorithm, respectively. To our knowledge, the analysis of a sequential version of the Mirror Prox algorithm has not appeared in the literature.